

* 05-PZQ-891 *





A. F. A. T. C. 10

EX 7

7

March - Los Angeles - December

CAT # 11-01

RENTAL PAGE COVER

ENGINEERING FORMULAE
THEIR MEANING AND DERIVATION



THIS BOOK IS PRODUCED
IN COMPLETE CONFORMITY WITH
THE AUTHORIZED ECONOMY STANDARDS

OTHER USEFUL TECHNICAL BOOKS

Spons' Electrical Pocket-Book. A reference book of General Electrical Information, Formulæ and Tables for Practical Engineers. By **W. H. Molesworth**, M.I.E.E., M.I.Mech.E., and **G. W. Stubbings**, B.Sc., F.Inst.P., A.M.I.E.E. **Ninth edition.** F'cap. 8vo. 410 pp., over 300 illus. 9s. net.

A Pocket-Book of Useful Formulæ and Memoranda for Civil and Mechanical Engineers. By **Sir G. L. Molesworth**, K.C.I.E., Past-President of the Institution of Civil Engineers, M.I.Mech.E. Edited by **A. P. Thurston**, D.Sc. (Eng.) Lond., M.I.Mech.E., M.I.A.E. **Thirty-second edition**, revised and enlarged. Oblong 32mo., 968 pp., over 750 illus. 12s. net.

Spons' Engineers' Tables for Civil and Mechanical Engineers, Builders, Contractors, Plumbers, etc. By **J. T. Hurst**, C.E., Mem. Phys. Soc. of London, late Surveyor War Department, etc., etc. **Thirteenth edition**, new impression, 288 pp., 64mo. (waistcoat pocket size). 2s. 6d. net. Leather, rounded corners, gilt edges. 3s. net.

Molesworth's Aeronautical Engineers' Pocket-Book. Comprising information and data useful to those engaged in the design of aircraft and their engines. By **A. P. Thurston**, D.Sc. (Eng.) Lond., M.I.Mech.E., F.R.Ae.S., M.I.A.E., and **W. O. Manning**, F.R.Ae.S. Ob. 32mo., 320 pp., 206 illus. 10s. net.

Barlow's Tables of Squares, Cubes, Square Roots, Cube Roots, and Reciprocals of all Integer Numbers up to 12,500. **Fourth edition**, revised and enlarged by **Dr. L. J. Comrie**, M.A., Ph.D., F.R.A.S. Demy 8vo., 272 pp. 10s. net.

Handbook of Formulæ, Tables, and Memoranda. For Architectural Surveyors, Draughtsmen, and others engaged in Building. By **I. T. Hurst**, Past-President Association of Surveyors H.M. Service. **Seventeenth edition.** New Impression. Oblong 32mo., 697 pp., 235 illus. 8s. 6d. net.

The Engineer's Sketch-book of Mechanical Movements, Devices, Appliances, Contrivances, and Details employed in the Design and Construction of Machinery for every purpose. Classified and arranged for reference for the use of Engineers, Mechanical Draughtsmen, Managers, Mechanics, Inventors, Patent Agents, and all engaged in the mechanical arts. By **T. W. Barber**, M.Inst.C.E. **Sixth edition.** Demy 8vo., 367 pp., 2987 illus. 12s. 6d. net.

The Automobile Engineers' Pocket Book of Rules, Tables, and Data. A textbook of proved formulæ for Designers and Students. By **H. Kerr Thomas**, M.I.Mech.E., M.I.A.E., etc. Fcap. 8vo., 126 pp., 50 illus. 7s. 6d. net.

ENGINEERING FORMULAE THEIR MEANING AND DERIVATION

AN ELEMENTARY TREATISE DEALING WITH PRINCIPLES
AND IDEAS, WRITTEN SPECIALLY FOR ENGINEERS TO
ENABLE THEM TO UNDERSTAND THE LANGUAGE
OF MATHEMATICS AND TO GRASP THE UNDER-
LYING MEANING OF ITS SYMBOLISM

BY

G. W. STUBBINGS

B.Sc., F.Inst.P., AM.I.E.E.



36890

5-26'58

6216 W. Va. Inst. of Tech.

LONDON

E. & F. N. SPON, LIMITED, 57 HAYMARKET, S.W.1

1943

QA
37
S88

PREFACE

THIS book has been written to explain in a simple manner the underlying principles of those branches of elementary mathematics that are used in engineering and applied science. Throughout the book, and particularly in the early chapters, the author has endeavoured to present the deductive reasoning of mathematical science in words, rather than in the highly-condensed symbolic form customary in the ordinary textbook. It is hoped that this method of treatment will be found helpful by many to whom the ordinary presentment of mathematics is difficult or even repellent, and who are thereby prevented from understanding something of what is probably the most interesting of all the sciences.


The main plan of this book is the development of the ideas of number and function. After preliminary chapters on the bases of algebra and geometry, the generalisation of the number idea to the symbolic representation of a directed step is developed. The idea of function is treated in a general way from the analytical and graphical standpoints, and the graphical treatment is made the basis of an approach to the elementary principles of the calculus. After a brief discussion of equations, periodic functions are studied at length. In the chapter on derived functions, a basic definition of the logarithmic function is laid down. This definition is shown to have a geometrical interpretation in reference to the properties of the rectangular hyperbola, and the properties of logarithmic, exponential and hyperbolic functions are derived from a study of the geometry of this curve. The final chapter deals briefly with the exponential function representing the normal law of errors, and touches upon the application of mathematics to statistics.

Whilst this book has been written primarily for those who wish to understand something about mathematics, and especially its practical applications to the theory of engineering, rather than to acquire facility in the technique of working out sums and solving problems, it should be helpful to students engaged on normal mathematical courses, as a supplement to their textbooks. The author hopes also that this book will be found useful by lecturers on elementary engineering mathematics.

G. W. S.

CONTENTS

CHAPTER		PAGE
I	INTRODUCTION	1
II	THE BASIS OF ARITHMETIC	8
III	GEOMETRY	30
IV	NUMBER	54
V	FUNCTIONS	85
VI	GRAPHS	106
VII	EQUATIONS	133
VIII	PERIODIC FUNCTIONS	156
IX	DERIVED FUNCTIONS	184
X	EXPONENTIAL AND LOGARITHMIC FUNCTIONS .	213
XI	THE APPLICATION OF MATHEMATICS TO STATISTICS	237
	INDEX	249



Digitized by the Internet Archive
in 2023

ERRATA

- Page 17, line 26. For " 39 " read " 37 ".
- Page 23, line 4. For " $1, 2$ and 3 " read " $1, 2$ and 1 ".
- Page 23, line 17. For " $\text{always to } 5$ " read " $\text{always to } 4$ ".
- Page 31, lines 27 *et seq.* Read " $\text{intersection of } c \text{ and } d,$
determines a point on a and a corresponding one on b .
Thus for every point on a there is one on b , so that,
although b is longer than a ,
- Page 34, line 28. For " b " read " a ".
- Page 35, line 30. For " restrictively " read " intuitively ".
- Page 43, line 16. For " $\frac{m^2}{n^4}$ " read " $\frac{m^2}{n^2}$ ".
- Page 45, line 9. For " $\frac{b}{c}$ " read " $\frac{c}{b}$ ".
- Page 63, line 22. For " $3\frac{113}{16}$ " read " $3\frac{116}{113}$ ".
- Page 67, last line. For " $10\frac{3}{10}$ " read " $10\frac{3}{10}$ ".
- Page 81, line 21. For " $x^2 + (jy)^2$ " read " $x^2 - (jy)^2$ ".
- Page 81, line 35. For " $\cos \theta$ " read " $\text{cis } \theta$ ".
- Page 82, line 36. For " $\cos^2 \theta$ " read " $\cos 2\theta$ ".
- Page 126, line 17. For " $\frac{x}{2a}$ " read " $\frac{2x}{a}$ ".
- Page 126, line 20. For " $\frac{x}{2a}$ " read " $\frac{2x}{a}$ ", and for " $2an$ "
read " $\frac{1}{2}an$ ".
- Page 126, line 22. For " $y = b + 2n^2a - 4n^2a = b - 2n^2a$ "
read " $b + \frac{1}{4}n^2a$ ".
- Page 136, line 24. For " $ax + by + c$ " read " $bx + ay + c$ ".
- Page 136, line 27. For " $a_1x + b_1y + c_1$ " read " $b_1x + a_1y + c_1$ ".
- Page 137, line 27. For " 1 " read " 2 ".
- Page 141, line 2. For " x^3 " read " x^2 ".
- Page 145, line 27. For " $x^2 - 2ab + b^2$ " read " $x^2 - 2ax + a^2$ ".
- Page 165, line 19. For " $(\cos^2 \theta)^2$ " read " $\cos^2 \theta$ ".
- Page 167, line 33. For " ϕ " read " θ ".
- Page 168, line 1. For " ϕ " read " θ ".
- Page 182, line 8. For " $b^2 \sin(p\omega t + \theta)$ " read " $b^2 \sin^2(p\omega t + \theta)$ ".
- Page 188, line 2. Delete " \sqrt{u} ".
- Page 195, line 25. For " ax^5 " read " ax^6 ".
- Page 220, lines 6 and 7. For " $x_1 + x_2$ " read " $\log x_1 + \log x_2$ ".
- Page 229, line 5. For " $e^{-j\theta}$ " read " $e^{j\theta}$ ".
- Page 229, line 21. For " e^{2j} " read " $e^{2\pi j}$ ".

ENGINEERING FORMULAE

THEIR MEANING AND DERIVATION

CHAPTER I

INTRODUCTION

Nature and Utility of Mathematics. Mathematics is commonly said to be the science of quantity, that is, the branch of knowledge that enables us to answer questions of the kinds, How many ? or How much ? Like other sciences, mathematics has its theoretical and applied sides. Theoretical or, as it is sometimes called, pure mathematics is a logically developed system of knowledge based upon common notions intuitively perceived, and upon postulates or ideas which are not self-evident or necessarily true, but which are taken for granted. A science of this kind is concerned with the idea of quantity in an abstract kind of way, that is, apart altogether from the sizes or numbers of any particular material things. Practical or applied mathematics is concerned with the application of pure mathematics to material phenomena ; what is called the theory of engineering is essentially a branch of applied mathematics. Pure mathematics has been considered by some thinkers to be a barren science, because of itself it yields no material results. According to this view the study of mathematics is meritorious only in its applications to material things, and mathematics is called the handmaid of the useful sciences. Other thinkers have seen in the study of pure mathematics an example of that unselfish pursuit of truth which is supposed to be one of the highest attributes of mankind.

Engineers and Mathematics. Engineers naturally are concerned with the applied or practical side of mathematics in so far as, professionally, they are concerned with this subject at all. The question, How much mathematics does an engineer require in this profession ? is one that has never been satis-

factorily answered ; indeed a complete answer is impossible because the needs of engineers are so diverse. Some engineers, such as are engaged in design or research work, require to know a large amount of mathematics ; others engaged in different branches of the profession get along and achieve success with a very slender mathematical knowledge. The courses of study in engineering schools seems to be designed with the idea of striking an average between these two extremes, with the consequence that most engineers learn, or rather are taught, much more mathematics than they seem to require in their later years. This superfluous knowledge which they acquire is generally rapidly forgotten, particularly so because the process of learning is accompanied with severe mental effort in the working out of sums the like of which never seem to crop up in after life. The engineering student juggles with algebraic and trigonometrical expressions, he solves equations, and often learns how to differentiate and integrate according to the rules of the calculus. In later practical life he has his sums largely worked out for him in books of reference, and his practical mathematical work is confined in the main to the use of a slide rule and the plotting of an occasional graph. It would appear at first sight, therefore, that much of the mathematics taught in the engineering schools is useless in the after life of the greater part of the students. This view, however, is not quite correct. Although relatively few engineers are called upon to solve difficult problems or to work out hard sums in the practice of their profession, they should be able to understand the implication of the mathematical parts of the technical literature they read, in the same way as they are able to understand the meaning of a drawing. This ability to understand what is meant by a mathematical argument, and to give mental assent to it, calls for mathematical knowledge of a kind, but not necessarily for the ability actually to solve problems and to work out hard sums. In short, as an engineer can understand the drawings of a machine without being able to design it, so he ought to be able to understand mathematical investigations of a relatively simple character without being able to originate them. This ability to understand what is meant by

mathematics and to give it rational mental assent is quite a different thing from ability to solve problems. In short, in the field of mathematics applied to engineering, we must carefully distinguish between principles and theory, and practice or technique. Let us consider briefly these two ideas, mathematical theory and mathematical technique.

Theory and Technique of Mathematics. A certain amount of mathematical technique is necessary for all but irresponsible people who live in a civilised country. If, for instance, a person buys two articles in a shop, costing say 2s. 2d. and 4s. 5½d., and he tenders a 10s. note in payment, he must be able very quickly to calculate his total indebtedness, 6s. 7½d., and the amount of change, 3s. 4½d., he should receive. Further, if this change is given to him as a half-crown, a sixpence, a threepenny piece, a penny, and a halfpenny, he should be able, at a glance, to check that the total is correct. This sort of technique is taught to children in elementary schools, together with a great deal more, such as the working of sums in long division, and with our complicated system of weights and measures this is considered by many educationists to impose upon them mental burdens heavier than they ought to bear.

Let us consider an example of this kind of elementary technique, that of long division. An intelligent child may ask why such and such is done in order to work the sum, and why the process gives a correct result. The answer to this question involves the other branch of mathematics—that of theory, and it is here important first to realise that a knowledge of the rational basis of the technical process is not necessary for working the sum. The majority of children who acquire the technique of working such sums correctly have no idea of the reasons for what they do. Secondly, it is possible to have a perfectly clear grasp of the underlying theory of long division, and yet to be deficient in a peculiar kind of ability that is needed to carry out a lengthy calculation accurately.

Mathematical technique of a higher order than that taught in elementary schools can be acquired without any firm grasp of underlying principles. Quadratic equations in algebra

can, for instance, be solved by the mechanical use of a formula, and many mathematical expressions can be differentiated by a mechanical application of the rules of the calculus. There is, however, a much higher technique, that of the use of mathematics as an instrument of research in the discovery of new knowledge. This technique depends to a large extent on an inborn talent, and a person possessing this talent can rightly be called a mathematician. The use of mathematics in this way calls for intuition which, although it is developed by training and practice, can never be taught unless the inborn talent is present.

The ultimate object of all practical or applied mathematics is the evaluation of numerical results. The final technique of mathematics is thus arithmetical, and consists in the working-out of ordinary sums. We have seen that the ability to do work of this kind is quite a different matter from an understanding of the processes required for the work ; it is also different from the highest technique of mathematics, that of research ; indeed many mathematicians of undoubted genius have confessed their inability to work long arithmetical sums correctly. A large amount of arithmetical calculation is now avoided by the use of the labours of bygone mathematicians who have spent years of their lives in working sums and tabulating the answers. Tables of squares and cubes, of logarithms, and of sines and tangents of angles are examples, which will be familiar to engineers. Further, a large amount of arithmetical technique is now carried out by machines. In the commercial field, columns of figures are recorded and automatically added by machines called comptometers, and more difficult arithmetical sums can be worked with great rapidity and absolute accuracy by various calculating machines. It is worthy of note, however, that the use of these latter machines requires a new technique, and a new species of scientific worker has arisen for the exercise of this technique, the professional computer.

Difficulties of Mathematics. We have seen that in the field of mathematics there are three distinct kinds of ability, first that of the technique of the application of rules of procedure to the working of sums or the solving of set problems, secondly,

that of the use of mathematics as an instrument of research in the discovery of new knowledge in applied science, and, thirdly, in that of understanding readily a mathematical argument in published form. This latter ability is of particular importance to engineers and practical scientists in order that they may be able to grasp and give mental assent to the original work of those who, having used mathematics as an instrument of research, publish the results of their labours for the benefit of practical men. The question has often been asked, Why is it that mathematical arguments are so difficult and even repellent to engineers who are supposed to have received a mathematical education in their student days, that they almost invariably skip in their reading when they encounter mathematical symbolism? There is little doubt that this may be to a small extent the fault of the original workers and writers themselves, in that they are not sufficiently careful to present their mathematics in the most lucid form, or even that they may make an unnecessary display of mathematical symbolism in order to parade their erudition. This, however, only touches the fringe of the question. A very real difficulty in following the most lucid argument presented in a mathematical way is experienced by most engineers, and the reason for this difficulty is not at all clear. A mathematical argument is of all arguments the most logical and orderly, the meanings of mathematical symbols, unlike those of many words, are perfectly definite and free from equivocation, and the argument is free from the irrelevancies that encumber those presented in a verbal form. It would seem on the face of it that of all arguments, one of a mathematical character should be the easiest to follow, This, of course, is not the case, and the reason for this seems to be two-fold. First, the closely-knit and extremely compact nature of a mathematical argument requires a considerable amount of mental concentration for its appreciation, and secondly, the symbolism in which the argument is expressed is really a special kind of language, and the implications of this symbolism must be so familiar that they are grasped instantly and subconsciously, without any effort of the memory. The importance of this latter point will be realised

by reflecting that a letter written in English but with Greek characters would be read with difficulty by anybody but a Greek scholar to whom the association of Greek and Roman letters is intuitive and without conscious effort of memory.

Mathematical Symbolism. The symbolism of mathematics is primarily a kind of universal language, the elements of which are familiar to most people. Thus, all but the illiterate know that the figures 3 and 2 in juxtaposition mean a combination of 3 tens added to 2 units. Further, nearly everybody knows that the meaning of the combination of the two figures can be changed by the use of a suitable mathematical symbol : thus $3 + 2$, $3 - 2$, 3×2 , $3 \cdot 2$, 32% , and 3^2 all have different meanings. Mathematical symbolism is, however, of a more subtle character than a mere shorthand, and it can be said that the development of mathematical science has been in a large measure due to the development of symbolism, particularly in the extension of the application of symbols beyond what may be called the common-sense field of their application. Thus the expression $3 - 2$ is perfectly clear to all literate people ; it is a short way of writing, subtract 2 from 3. But $2 - 3$ is inherently meaningless. Similarly 3^2 is a short way of writing, multiply two 3's together, but an expression like $3^{\frac{1}{2}}$ is also inherently meaningless. The development of mathematics has been largely due to the discovery, or rather the assignment, of meanings to operations which, from a common-sense point of view, are absurd. This does not mean that the development of mathematics is anarchical ; a meaning for an apparently absurd expression like $2 - 3$ is only adopted when it is found to lead to results consistent with the well-established branches of the science. Thus, as we have already stated, mathematical symbols are something more than mere shorthand, and they are often the signs that stand for profound ideas. A proper realisation of this subtle character of mathematical symbols, and of the extension of their meanings beyond the elementary and common-sense, is one of the essentials for the ability to understand mathematical theory, and, hence, to follow and give mental assent to a mathematical argument.

Object of this Book. Having set down the foregoing intro-

ductory remarks about the nature of mathematics we can outline the object of this book and the plan according to which it has been written. Its object will not be to teach any technique or to provide hints or practice for the working of sums or the solution of mathematical problems, but rather to explain mathematical theory and the implications of mathematical symbolism in order to assist the engineer to understand and read the universal language of applied science and to give mental assent to what this language is intended to convey. We shall see that the greater part of the mathematical theory required by engineers for this purpose can be studied under two main heads, the extension of the idea of number from its primary significance in association with the basic operation of counting, and the conception and nature of functional relationship or of the quantitative and mathematical expression of the law of causality, that is, of the assumption upon which the whole art of engineering is based, that the same causes always produce the same effect.

CHAPTER II

THE BASIS OF ARITHMETIC

Counting. The simplest and most fundamental of all mathematical ideas is that of counting, and one of the earliest landmarks in the education of a young child is his ability to count up to twenty. This idea of counting is so firmly implanted in the minds of most people that its inherent nature is rarely perceived clearly. When we say a child can count up to twenty, what do we mean? Simply this, that the child has learned by rote and can correctly repeat the words one, two, and so on, to twenty, in their correct order; it is not sufficient for the words to be memorised, the order is all-important. The child uses his knowledge of this sequence of words to count a collection of say eleven objects by touching each in turn and repeating the words in order, one for each object. The word repeated when the last object is touched is the number of the collection. Otherwise, the words might be written down, and the objects put in what is called one-to-one correspondence with the words, starting at the first word of the series. The last word associated with an object would then give the number of the collection. The actual words in the ordered sequence are more or less arbitrary, it is their invariable order that matters. The number words in the French language are, for instance, different from those in English. There is no end to this sequence of cardinal or natural number words, as they are called, because, however many objects there may be in a collection we can always conceive that there might be more. If the words are written in the usual direction, left to right, each word is characterised by its left hand neighbour, which it follows, with one exception, the starting word, one in English, which being the origin has no left hand neighbour. The number words are indicated by special symbols called figures or digits.

Numeration. If number words and symbols were arbitrarily chosen, as they might conceivably be, the vocabulary required for counting would be unlimited in extent. Economy in the

use of number words and symbols is obtained by the nearly universal custom of counting in tens. According to this system, to count a collection we collect it into sets of ten, and count the residue, if any, less than ten. Sets of ten are then collected into larger sets of ten tens, and so on. For this method of counting we require the number symbols for the words one to nine inclusive and also an additional symbol to represent none or nothing. We can then write ten in the usual way 10, which by a convention of marvellous utility, means one set of ten and nothing over. Similarly 101 means 1 large set of ten tens, no smaller sets of ten, and one over. Numbers symbolically expressed in this way are read from right to left, and the starting number denotes odd singles or units, the next tens, the next tens of tens or hundreds, and so on. This extremely convenient system of numeration, it should be noted, depends upon the use of a special symbol 0, called zero, to denote nothing.

It is generally supposed that the system of reckoning by tens is a development of finger-counting. Although this system is so firmly established it is not the best that could be devised, because the number ten can only be equally divided in two ways, into five two's and two fives, and of these, division into five parts is rarely wanted in every-day life. A system of counting by twelves would really be much more useful than that of counting by tens, because a collection of twelve objects can be divided into so many equal parts, six, four, three, and two. The system of counting by tens is called the decimal system, that of counting by twelves, the duodecimal system. It is worth noting that a duodecimal system of notation would require special symbols for ten, the number following nine, and for eleven, the number following ten. If these symbols were respectively t and e , then 10 would mean one set of twelve and nothing over or one dozen, $1t$ would mean a dozen and ten, and $1e$, a dozen and eleven, or the number indicated in the decimal system as 23. The method of reckoning by dozens and dozens of dozens or grosses is used commercially to some extent, but as there are no proper symbols for ten and eleven this is not a system of numeration. Another possible and rather interesting system of numeration is the binary, that

based on reckoning by pairs. In this system two number symbols only would be required, 1 and 0. 10 would stand for a pair, the number represented in the decimal system by 2, 11 would be 3 in the decimal system, a pair and one over, and 100 would be a pair of pairs or 4, decimal. The apparent simplicity of this system is offset by the very large number of symbols or digits required to represent quite small numbers, thus 33 decimal would be 100001 binary.

In nearly all civilised countries measures of concrete quantities such as length, bulk, weight, and money are all based on a decimal system. In the British system of weights and measures we reckon by 12's, 3's, and $4\frac{1}{2}$'s in length measures, in 16's, 28's and so forth in weight measures, and in 4's, 12's, and 20's in money. The utility of a decimal system for these kinds of reckonings is manifest, and it has been talked about in the country for the past hundred years, without result. The Englishman is so proud of his supposed capability to muddle through, that it seems almost as if he takes a pride in muddle as muddle.

Addition. Next to counting, the most elementary idea of mathematics is that of adding. The fundamental problem of addition is this: given the numbers of two collections, what is the number of the larger collection formed by combining them? The operation of solving this problem is simple. In the sequence of natural numbers we take as a starting point the number next after that of one collection and, from this start, we proceed along the sequence by steps corresponding to the number in the other collection. The number we reach finally is that of the combination, or the sum of the numbers in the component collections. It is necessary to note here an important point: it does not matter with which component number we start, the sum will always be the same. Thus if we represent one number by a sign or symbol a , and the other by b , and the operation of addition by the usual sign $+$, then $a+b$ is the same as $b+a$; otherwise, adding b to a gives the same result as adding a to b . We express this even more concisely by the shorthand method $a+b=b+a$. Let us consider the number obtained by adding b to a , and think of this being added to a third number represented by the sign c .

We denote the compound number which is the sum of a and b by the symbol $(a+b)$, the bracket showing that we are thinking not of two separate numbers but of the single number which is their sum. The operation of adding the number $(a+b)$ to c can be represented $c+(a+b)$, and we know, or can easily check by simple trials, that the result of this is the same as that represented by $c+a+b$ which means add a to c and to the answer add b . In shorthand, or what are generally called algebraic symbols, we can therefore write:—

$$c+(a+b)=c+a+b,$$

This all may seem rather trite, but trite things are often worth a little thought. The two rules of addition considered are, first, that the order is indifferent, and secondly, that in adding a number which is itself a sum, we can either add the whole sum at once, or we can add the two components of this sum separately. In technical language this first rule is expressed by saying that addition obeys the Law of Commutation, the second, by saying that it obeys the Law of Association. We shall see in the course of our study that a right understanding of mathematics depends to a considerable extent on a clear recognition of when these laws are or are not obeyed.

There are two further points about addition to be noted. However large a number is we can always conceive that same number can be added to it to make it larger. In other words the range or field of addition is unlimited or unrestricted. Symbolically, the shorthand expression $a+b$ always has a meaning.

The second point to be noted is that in addition the symbol zero or nothing, 0, has the properties of an ordinary number. Adding nothing to a number leaves it unchanged, so symbolically $a+0=a$, and also $0+a=a$, that is, adding a to nothing gives an answer a . This also seems trite, but a clear understanding of when 0 can and cannot be treated as an ordinary number is also very important in the study of mathematics.

Subtraction. Subtraction is inherently the reverse of addition. To subtract a number b from another number a consists in finding a third number c which when added to b will give a sum equal to a . To subtract b from a we count backwards in the sequence of natural numbers starting at a ,

by steps equal in number to b . We note first that we only obtain a result if a is greater than b , in short we cannot subtract from a number of objects a larger number. The operation of subtraction is restricted, $a - b$ only gives a numerical result if a is greater than b . Suppose, however, that b is a ; the operation of subtraction then leaves nothing, so that $a - a = 0$. Subtraction, to give a result at present intelligible, is therefore restricted by the condition that b must not be greater than a .

We see at once that subtraction, unlike addition, does not obey the commutative law. $3 - 2$ is essentially different from $2 - 3$. A number which is the sum of two numbers a and b can, however, be subtracted from a third number at one step, or in two steps, provided that $(a + b)$ is not greater than c . In other words, subject to this condition, $c - (a + b) = c - a - b$. Here the law of association is obeyed.

Let us now think of a number which is the result of subtracting one number from another. We can represent the compound number by the symbol $(a - b)$, the brackets here indicating as before that we are considering a single number. Now we can subtract the number $(a - b)$ from a third number c provided that c is the greater of the two. Assuming this condition is satisfied the second subtractive operation can be indicated symbolically as $c - (a - b)$. Is there any law of association that holds here? In the case of addition, and in the case of compound subtraction in the preceding paragraph, the law of association meant simply this, that the symbolic expressions $c + (a + b)$ and $c - (a + b)$ could be converted by removing the brackets and affixing to each number inside these brackets the sign $+$ or $-$ attached to the compound number signified by these brackets. In the first case the component numbers a and b are successively added to c , in the second case they are successively subtracted. The meaning of $c - (a - b)$ is not, however, at first sight clear, but this meaning can easily be discovered by experiment. Let us consider a number $(4 - 3)$ subtracted from 8. We know that as $(4 - 3) = 1$ the answer must be 7. If we write $8 - (4 - 3)$ and work the sum in two stages we see at once that first taking 4 from 8 gives less than the correct answer because we have subtracted too

much. How much too much? Evidently 3, because 4 is greater than $(4 - 3)$ by 3. Hence, after subtracting 4 from 8 and obtaining the first answer 4 we must add 3 to get the correct final answer 1. Symbolically $8 - (4 - 3) = 8 - 4 + 3$. This is quite general, and can be expressed generally in the shorthand or algebraic method as $c - (a - b) = c - a + b$. To remove the brackets in $c - (a - b)$ we must therefore alter the - sign of b to +, and it is only subject to this condition being fulfilled that the law of association is satisfied. The operation of subtraction has evidently led us into deeper water than we found when considering that of addition.

Multiplication. Multiplication is essentially repeated addition; to find the number in say four sets of five objects we count the total number in the four sets of five. This operation is denoted symbolically by 5×4 when figures are concerned; when we are considering multiplication in a general way, and we indicate any pair of numbers by letters a and b , the corresponding symbolism $a \times b$ is generally simplified by writing the product as ab . This is one of the inconsistencies of mathematical symbolism which should be noted carefully, ab means $a \times b$, but if definite values, say 4 and 3, are respectively assigned to a and b , ab cannot be written 43; we must use the full form with the \times sign and write 4×3 , because 43 means 4 tens + 3.

Let us now consider the operation of multiplication in the same way as we have considered those of addition and subtraction. In the first place we see at once that its field is unrestricted—since it is merely repeated addition. Secondly, multiplication obeys the law of commutation, in other words a multiplied by b is equal to b multiplied by a . This, although not perhaps self evident, is readily checked by experiment, and is confirmed by the fact that in reckoning the size of a rectangle, which depends upon the product of length by breadth, it is indifferent which side we take as the length.

Repeated multiplication obeys the law of association, for if a number c is equal to ab , then cd is equal to $c \times (ab)$ and to $c \times a \times b$. Multiplication also obeys the associative law when combined with addition, for if a number c is multiplied by a number equal to the sum of a and b , then writing this operation

as $c(a+b)$ this is equal to $ca+cb$. ca and cb are called the partial products in the multiplication. This associative property of multiplication is the basis of the ordinary method of multiplying a number by one less than 10, or of short multiplication. The thousands, hundreds, tens, and units of the number are multiplied separately and the partial products are added together mentally and written down as they are found.

We can usefully consider the associative properties in relation to the more complicated idea of the multiplication of two sums. Suppose these sums are $(a+b)$ and $(c+d)$. Then first treating $(a+b)$ as a single number, we have

$$(a+b) \times (c+d) = (a+b) \times c + (a+b) \times d = c \times (a+b) + d \times (a+b)$$

and, by obtaining the secondary set of partial products, we obtain finally

$$(a+b)(c+d) = ac + ad + bc + bd.$$

If the two numbers to be multiplied together are equal, say $(a+b)$, we find

$$(a+b)(a+b) = a \times a + ab + ba + b \times b.$$

A number a multiplied by itself is denoted by the symbol a^2 called " a square." Further, as the two partial products ab and ba are by the commutative law equal, we can lump them together in one number or term and write their sum as $2ab$. Thus we find

$$(a+b)^2 = a^2 + 2ab + b^2.$$

Let us now consider multiplication in relation to subtraction as well as addition. We have already seen that we can add a number $(a+b)$, itself a sum, in one or two stages, or that $c+(a+b)=c+a+b$. If the number to be added is itself a difference like $(a-b)$ we can find the result of adding this to c by first adding a and then subtracting b , otherwise $c+(a-b)=c+a-b$. Let us now consider an interesting case of the multiplication of a sum and a difference of the same numbers.

$$(a-b)(a+b) = (a-b) \times a + (a-b) \times b.$$

This is equal to $a \times (a-b) + b \times (a-b) = a^2 - ab + ba - b^2$.

In working out this final sum we have to subtract ab and add ba , and as these two products are equal, the effect of these two operations is simply zero or nothing. We therefore have, finally, $(a+b)(a-b)=a^2-b^2$. This result, which shows that the difference of the squares of two numbers is equal to the product of their sum and their difference, is of great importance in mathematics. The serious reader ought to test it by working an example in numbers: for example $9 \times 9 - 7 \times 7 = 81 - 49 = 32$, ought to equal $(9+7) \times (9-7)$ or 16×2 , as it evidently does.

We have seen that the compound subtraction of subtracting a difference indicated by $c-(a-b)$ must be interpreted as $c-a+b$. Similarly the operation $c-d \times (a-b)$ has the meaning $c-da+db$, provided of course that $d \times (a-b)$ is not greater than c .

Finally, let us consider the nature of 0 or zero in relation to multiplication. What is the meaning of $a \times 0$? It means taking nothing a times, and is evidently the same as 0 itself. So also is $0 \times a$ or taking a no times. Thus $a \times 0 = 0 \times a = 0$, whatever a may be. It is evident that in multiplication, zero or 0 is not an ordinary kind of number, because multiplying 0 always gives the same answer, 0.

We see from the foregoing that, although multiplication, being merely repeated addition, is inherently simple, yet the complete implications of this operation leads to some rather complicated results that require serious consideration for their firm mental grasp.

Division. Division is essentially the reverse of multiplication. To divide a number a by a second number b , we have to find a third number c , such that $c \times b = a$. Division is indicated by the symbol \div , $a \div b = c$ means therefore that

$a = c \times b$. Division of a by b is also denoted by the symbol $\frac{a}{b}$,

but as we shall see later, this symbol is further reaching in its implication than the mere operation of division in the rudimentary sense of the term. Division can be considered also as repeated subtraction. Thus, to divide a by b , we find the number of times b must be subtracted from a till nothing is left. We see at once that the field of division is restricted.

We cannot divide a number by one greater than itself in the basic sense of division, although later we shall see what meaning can be given to this idea. Further, a number cannot always be divided exactly by one less than itself; subtracting the dividing number successively we shall generally find that when the last possible subtraction has been made something is left over. This is called a remainder. As a matter of fact it is evident that a number has but a limited number of exact divisors, and some numbers have none at all. We shall return to this point later.

The operation of division is not commutative for, plainly, $a \div b$ is not the same as $b \div a$. In relation to addition, moreover, the associative law is only partly obeyed. Thus, while we can divide a sum $(a + b)$ or a difference $(a - b)$ by a number c in two stages, so that $\frac{a+b}{c} = \frac{a}{c} + \frac{b}{c}$, provided of course that a and b can both be divided by c , we cannot say that c divided by $(a + b)$ is equal to $\frac{c}{a} + \frac{c}{b}$. Thus $\frac{(6+4)}{2} = \frac{10}{2} = 5$
 $= \frac{6}{2} + \frac{4}{2} = 3 + 2$, but $\frac{12}{(4+2)} = \frac{12}{6} = 2$ is not equal to $\frac{12}{4} + \frac{12}{2} = 3 + 6$.

If a division of a sum is carried out in two stages, as $\frac{a+b}{c}$
 $= \frac{a}{c} + \frac{b}{c}$, $\frac{a}{c}$ and $\frac{b}{c}$ are called the partial divisions or partial fractions.

Let us now consider whether zero or 0 behaves as an ordinary number in relation to division. What is the meaning of $a \div 0$? Suppose there is an answer, b , then by the fundamental definition of division we must have $b \times 0 = a$. But we have seen that all numbers, however large, when multiplied by 0 give the same answer 0. There is, consequently, no answer to the sum $a \div 0$, and zero cannot be considered as an ordinary number in connection with division. This point is of very great importance, should be firmly grasped, and should be compared with the relations of 0 to the other basic arithmetical operations. Zero enters into addition and subtraction as an ordinary number, in multiplication zero is a number giving an answer which is definite, but anomalous.

Division by zero, however, is impossible because no number exists that will express an answer. In reference to division, therefore, zero is not an ordinary number. Sometimes a statement is made that division by zero gives an answer infinitely large, and this statement is expressed symbolically as $a \div 0 = \infty$, ∞ being the sign for infinity. A statement of this kind cannot, however, be entertained at this stage, although later in our study we shall see that, if the meaning of 0 is modified, the symbolic equation $a \div 0 = \infty$ may have some kind of significance. So long as zero is taken as meaning absolutely nothing, the result for instance of taking a from a , the operation of dividing by 0 must be considered as meaningless and impossible.

Factors, Factorials, and Primes. If a number a is the product of two numbers b and c , then b and c are called the factors of a . Given the factors of a number, the number itself can be found very easily. The reverse process of finding factors is, however, not so simple. We have already encountered some examples of how the symbolic or algebraic method of representing compound numbers enables us to find factors, thus, as $(a^2 - b^2) = (a + b)(a - b)$, it follows that if a number is the difference of two squares we can find two factors of it, the sum and difference of the numbers squared. Thus 1591 is equal to $40 \times 40 - 3 \times 3$, so that the factors of 1591 are $(40 + 3)$ and $(40 - 3)$. The serious reader will check this by multiplying 43 by 39.

A number formed by multiplying the natural numbers, or integers as they are often called, in sequence is called the factorial of the last number of the sequence. Thus $1 \times 2 \times 3 \times 4$ is factorial 4, and is written $4!$. Factorials are of great importance in mathematical theory, and the reader should make himself thoroughly acquainted with the factorial notation.

A natural number which cannot be divided into equal parts, or which has no factors, is called a prime. Strictly speaking, every number has two factors, itself and 1. Thus 7 can be taken as 1 part, 7, or can be divided by 7 to produce 7 ones. The better definition of a prime therefore is a number which has no factors but itself and unity or one. Thus 1, 2, 3, 5, and 7 are the primes less than 10. The further we proceed

along the sequence of natural numbers the less frequently do we encounter a prime number. Thus, between 90 and 100 there are only two primes, 93 and 97. We could evidently find the primes between 1 and any other number by striking out of the number sequence every second, every third, every fourth and so on. As the primes get fewer and fewer, the interesting question arises: Is there a last prime? This question can easily be answered by using a little algebraic symbolism. Suppose we find the last prime n . Form a number by adding factorial n to 1, that is $(n! + 1)$. If this be divided by any number less than n , there must be a remainder of at least 1. The number $(n! + 1)$ is therefore either a prime or is divisible by some prime greater than n , so that however large a prime number may be there must be a greater one. This was proved in the 9th book of Euclid's "Elements." A number which can be resolved into factors other than itself and unity is said to be composite. If two numbers are resolved into their smallest factors, and the factors, other than 1, of one number are all different from those of the other, the two numbers are said to be prime to each other. Thus 12 is prime to 77, but neither 12 or 77 are prime numbers.

Divisibility of Numbers. In practical calculations it is often desirable rapidly to determine whether a number is exactly divisible by a small number less than 10. Everybody knows that numbers divisible exactly by 2 are called even, and all other numbers are called odd. The general symbol for an even number is $2n$ because it is the double of some other number; the general symbol for an odd number is $2n + 1$ because it is 1 greater than some even number $2n$. All even numbers end in 2, 4, 6, 8, or 0. Any exact number of 100's is exactly divisible by 4, and every number greater than 100 is an exact number of hundreds plus a number represented by its last two figures. If this latter number is divisible by 4, the whole number is. Thus 47732 is divisible by 4 because 32 is, but 47754 is not exactly divisible by 4 because 54 is not.

A number is exactly divisible by 9 if the sum of the digits or figures expressing it is so divisible. This statement is so useful in practical calculation, and its proof affords so good an example of the utility of symbolic or algebraic calculations

that we shall devote a little space to it. Consider any number represented by three figures. a , b , and c , in order, where a , b , and c are figures or digits. The value of this number is $100a + 10b + c$, because the right-hand figure c represents ones or units, the next tens, and so on. Now $100a$ is equal to $99a + a$, and $10b$ is equal to $9b + b$. The whole number is therefore equal to $99a + 9b + a + b + c$. Now the part $99a + 9b$ is exactly divisible by 9, for 99 and 9 both are. The whole number will therefore be divisible by 9 if the second part $a + b + c$ is. But this second part is simply the sum of the figures or digits comprising the number. A further important point is that the remainder obtained when the number is divided by 9 is the same as the remainder obtained by dividing the sum of the digits, $a + b + c$, by 9. For example, consider the number 472185. The sum of the digits $4 + 7 + 2 + 1 + 8 + 5$ is 27, and this is divisible by 9. The number is therefore so divisible exactly. Again, consider 3538. The sum of 3, 5, 3, and 8 is 19, the sum of these digits $9 + 1$ is 10. The sum of 1 and 0 is 1. This is therefore the remainder when 3538 is divided by 9. These statements should be tested by actual division. The operation of finding in this way the remainder when a number is divided by 9 is called "casting out the nines," and it is the basis of a very old set of rules for checking the accuracy of the working of arithmetical sums.

This latter point is worth illustrating. Every number is the sum of so many nines plus the remainder obtained when it is divided by nine. Suppose two numbers are multiplied together, these can be expressed symbolically as $(9a + b)$ and $(9c + d)$, where b and d are the remainders when the two numbers are divided by nine. Let us multiply these symbolic expressions together. We multiply $(9c + d)$ first by $9a$, obtaining $81ac + 9ad$, and then by b , obtaining $9bc + bd$. The sum of these two partial products is $81ac + 9ad + 9bc + bd$. The sum of the first three terms is exactly divisible by 9, and the remainder when the whole product is divided by 9 is the remainder when bd is so divided. Hence, if a multiplication sum is correct the product of the remainders obtained by casting out the nines in the numbers multiplied, b and d in our example, is equal to the remainder obtained when the

nines are cast out of the product. For example, $384 \times 42 = 16128$. Cast out the nines in 384 and we obtain $3 + 8 + 4 = 15$, $5 + 1 = 6$. Do the same to 42: $4 + 2 = 6$. Multiplying these results we get $6 \times 6 = 36$, and $6 + 3 = 9$. Now cast out the nines in the product 16128: $1 + 6 + 1 + 2 + 8 = 18$ and $8 + 1 = 9$. This does not prove the correctness of the answer, for if we had made a mistake of any number of nines in the working the test would be satisfied. Thus an inaccurate answer 15228 would give a remainder 9. If, however, the check by casting out nines is satisfied it shows that the answer is probably correct.

The curious property of the number nine has nothing to do with the number itself. It is simply due to our system of reckoning by tens. If we reckoned by twelves and used the duodecimal system, we should, to check sums, cast out the number one less than twelve, that is the elevens or *e*'s.

Involution. We have already seen that a^2 is a short way of writing $a \times a$. The process of multiplying a number by itself is called evolution, and it evidently can be carried on indefinitely. The small 2 in a^2 is called an index, and it denotes the number of a 's that have to be written down with \times signs between them. Similarly a^3 means $a \times a \times a$, and generally a^n means n a 's set down with \times signs between them. One or two points about this idea of involution may be noted at the outset. First, however many 1's we multiply together the answer is always 1. Symbolically $1^n = 1$ where n is any number however large. Secondly, as a^2 denotes setting down two a 's and multiplying them, we can consider that a^1 means setting down one a only, and, in short, is simply a . Otherwise $a^1 = a$. Thirdly, as multiplying 0's together can only give answer 0, the operation 0^n , as far as it has any meaning at all, can only be considered to be equivalent to 0. The third power a^3 of a number is usually read " a cube." a^4 is read " a to the fourth," and generally, a^n , " a to the n th."

If two powers, say a^2 and a^3 of the same number a , are multiplied together, this means $(a \times a) \times (a \times a \times a)$ which is plainly the same as $a \times a \times a \times a \times a$, or a^5 . Hence the multiplication of two powers of the same number gives a third power the index of which is the sum of the indices of the numbers. Symbolically $a^m \times a^n = a^{m+n}$. Again to divide one power, say

a^4 , by a smaller power, say a^2 , means divide $a \times a \times a \times a$ by $a \times a$. Now $a \times a \times a \times a$ is evidently $a \times a$ times $a \times a$, and to obtain the answer in dividing by a^2 we simply strike out 2 a 's of the repeated product $a \times a \times a \times a$. Generally the division of one power by a smaller power of the same number gives a third power the index of which is the difference of the indices of the numbers. Algebraically $\frac{a^m}{a^n} = a^{m-n}$. Division of powers at present only has a meaning when m is greater than n .

Suppose m is 1 more than n . The answer to the sum $\frac{a^m}{a^n}$ is $a^{m-n} = a^1 = a$. This is correct, because if we strike out all but 1 factor in the multiplication of n a 's we shall be left with a single a as the answer. Suppose we divide a^m by itself, a^m . Any number divided by itself gives an answer 1. But, according to the index rule for the division of powers $\frac{a^m}{a^m}$ ought to equal $a^{m-m} = a^0$. But a^0 is meaningless as it stands, it indicates put down no a 's and multiply them together. We can, however, give a^0 a meaning by saying that it equals 1, because this meaning is consistent with the rule for the subtraction of indices in the division of the one power of a number by another. This is a good example of what we shall find quite common in mathematics, the assignment of a meaning to a symbol which in its primary significance is unintelligible. The zero power of any number, however large, is therefore equal to 1; $a^0 = 1$.

Let us consider the idea of a power of a power, say $(2^2)^3$, that is the cube of 2 square, or the cube of 4. This is $4 \times 4 \times 4$ which is equal to six 2's multiplied together. It is easy to see that, generally to raise power to a power, we multiply the index of the power by that of the power to which it is to be raised. Symbolically $(a^m)^n = a^{m \times n}$.

We will now return to a more detailed consideration of an algebraic sum worked out early in this chapter, that of finding the value of the square of the sum of two numbers a and b . The result we obtained can be represented :

$$(a + b)^2 = (a + b) \times (a + b) = a^2 + 2ab + b^2.$$

An algebraic expression like $a + b$, which has two separate numbers added together or subtracted the one from the other, is called a binomial. The answer to the sum $a^2 + 2ab + b^2$ is called the expansion of $(a + b)^2$ because it gives its equivalent in a longer or more expanded form. The 2 attached to the term ab in the expansion means that, in working the sum, we obtained 2 partial products ab exactly similar. A figure of this kind is called a numerical coefficient or, more briefly, a coefficient.

Let us examine $a^2 + 2ab + b^2$ more closely. We see that it contains 3 terms, one more than the index of the power of $(a + b)^2$. The power of a in the first term is 2, in the second term it is 1. Can a be said to have a power in the third term? It can, for $b^2 = 1 \times b^2$, and as we have seen $1 = a^0$ so the power of a in the third term is 0. The powers of a decrease from left to right, and the powers of b increase. This is shown by writing the answer in this form, $a^2b^0 + 2a^1b^1 + a^0b^2$, which is cumbersome, but which shows how the powers of one number a in the binomial decrease and those of the other is increased from left to right. Further, it shows that the sum of the indices of a and b in each term is the same, 2. A compound algebraic expression about two or more numbers like a and b which satisfies the condition that the sum of the indices of powers in each term is the same, is said to be homogenous.

We can find the cube of $(a + b)$, $(a + b)^3$, or $(a^2 + 2ab + b^2) \times (a + b)$, by multiplying $a^2 + 2ab + b^2$ first by a then by b and adding the results. The first multiplication gives $a^3 + 2a^2b + ab^2$, the second $a^2b + 2ab^2 + b^3$. In all the partial products obtained in finding $(a + b)^3$ there are therefore 3 terms a^2b and 3 terms ab^2 , and the result is

$$(a + b)^3 = a^3 + 3a^2b + 3ab^2 + b^4.$$

We note here again that the expansion gives 4 kinds of terms, a number 1 greater than the index of $(a + b)^3$. Further the expansion is homogenous, the sum of the indices of each term is 3, the same as the index of the power to which $(a + b)$ is raised. Also the powers of a successively diminish and those of b successively increase from left to right.

Let us now consider an interesting point about numerical

coefficients. Take the original binomial $(a+b)$. This is the same as $(a+b)^1$. There are 2 terms, and as $a=1 \times a$, we can say the coefficients are 1 and 1. Similarly we can say the coefficients in the expansion of $(a+b)^2$ are 1, 2, and 3. Let us consider the following schematic statement :

$(a+b)^1$	Coefficients	1	1		
$(a+b)^2$	„	1	2	1	
$(a+b)^3$	„	1	3	3	1

We see quite easily that these coefficients follow a rule or law. Those in the second line are formed by placing between those in the first an additional one equal to their sum. Similarly those in the third line are formed by adding the first two, $1+2$, and placing 3 under and between them, and similarly adding the next pair $2+1$ and placing the 3 under and between them. According to this law the coefficients of $(a+b)^4$ ought to be 1, 4, 6, 4, 1, and the indices of the powers of the 5 terms ought to add up always to 5. Thus $(a+b)^4$ ought to equal :

$$a^4 + 4a^3b + 6a^2b^2 + 4ab^3 + b^4.$$

If the reader has the patience to do so, he can check the correctness of the statement by actually multiplying the expansion of $(a+b)^3$ by $(a+b)$.

This short discussion contains the germ of what is called the Binomial Theorem, which was first enumerated by Newton, and which is of far-reaching importance in pure mathematics. There is no need, at this stage, to pursue this discussion, excepting to draw attention to an important point. The first term in the expansion of any power, say, the n th of $(a+b)$, is a^n . This is evident. The next term has a numerical coefficient equal to the index of the power, the index of the power of a is one less than n , and the index of the power of b is 1. The second term is therefore $n \times a^{n-1} \times b$. The rules for finding the other terms are given in books on algebra, but these are not very important to the engineer.

The Order of a Number. The digits or figures used to express a number in the decimal system of notation give the numbers of the various powers of ten in which it is divided. The first figure from the right, indicating ones or units, gives

the number of the zero powers of 10, since $10^0=1$. The second figure gives unit powers of 10 since $10^1=10$. The third figure, representing hundreds, gives the number of the second power of 10 as $10^2=100$. The index of the highest power of ten contained in a number is one less than the number of digits or figures that it contains, and this index is an indication of the order of the magnitude of the number. Thus one million written down as 1,000,000 is equal to 10^6 and is of the sixth order. This index notation is a very convenient one for representing very large numbers in an approximate kind of way. For instance, the distance from the earth to the sun is, on average, a little over 90,000,000 miles. This number can be written as being about 9×10^7 miles. Astronomers who claim to have measured the size of the universe in a four-dimensional world may tell us that this size may be expressed in miles by a number like 10^{60} , a short way of expressing a number normally written as a 1 followed by 60 zeros.

Numbers expressed roughly like 9×10^7 can be rapidly multiplied by the index rule. Thus 9×10^7 multiplied by 6×10^9 would be $54 \times 10^7 \times 10^9$, which by adding the indices of the tens is 54×10^{16} . The order of a product, it will be seen, is at least equal to the sum of the orders of the numbers multiplied. It may be more. For instance, the product 54×10^{16} means 54 followed by 16 zeros, it will require 18 figures to express it in the usual way, so that its order will be 17, one more than the sum of the indices of the tens. Most engineers know that the order of a number is measured in an exact kind of way by what is called its logarithm. Logarithms are artificial numbers which have the property that the logarithm of a product is the sum of the logarithms of the numbers multiplied. According to the ideas of number which are the basis of this chapter it is easy to understand that the index of a power of ten is its logarithm, for 1000×100 or $10^3 \times 10^2$ is 100,000 or $10^{3+2}=10^5$, but the nature of the logarithm of a number, say between 10 and 100, involves considerations and ideas of number much more subtle and far-reaching than those which form the basis of this chapter. We shall return to this important subject hereafter.

Evolution. Evolution is the reverse of involution; it is repeated division, and the problem of evolution consists in finding a number of which a given number is a given power. The answer to a problem of this kind is called a root. The evolution process of finding a number which when squared will equal a given number is said to be that of finding the square root, similarly the cube root of a number when cubed will equal the number. Thus the square root of 4, written $\sqrt{4}$, is 2, because $2^2 = 4$. Similarly the cube root of 27, written $\sqrt[3]{27}$, is 3, because 3^3 or $3 \times 3 \times 3 = 27$. The sign $\sqrt{}$ is called a radicle and the small figure appended for all but square roots is the order of the root. Thus $\sqrt[5]{}$ means the fifth root.

The operation of evolution is, inherently, very restricted; we can only at present find the square root of numbers which are themselves squares. Thus, of the numbers between 1 and 20, only 1, 4, 9, and 16 have square roots, 1, 2, 3, and 4 respectively. These numbers, called square numbers, get rarer and rarer the farther we travel along the series of natural numbers for if n^2 is one square number, the next is $(n+1)^2$ and this, by the rule for finding the square of a binomial, is $n^2 + 2n + 1$. Thus the next square number to n^2 is formed by adding $2n+1$ to it. 1,000,000 is 1000 squared, so the next square to 1,000,000 is obtained by adding 2001. The distance between neighbouring square numbers therefore continually increases, and yet, as every number must have a square, there must be as many square numbers as ordinary numbers. The reader may well think about, but not worry about, this paradox.

As in the involution of powers we multiply indices, so in the evolution of powers we divide indices. Thus as $(a^m)^n$ is equal to a^{mn} , so the n th root of a^{mn} or $\sqrt[n]{a^{mn}}$ is simply a^m . This method of evolution is, at present, only intelligible when the index of the power to be evolved is exactly divisible by the order of the root. $\sqrt{2^4}$ or the square root of 16 is evidently 2 raised to a power which is equal to 4 divided by 2. We can at this stage, however, assign no meaning to the idea of, say, expressing the cube root of a^4 by dividing the index 4 by the order of the root, 3.

Algebraic Symbolism. We have, in this chapter dealing primarily with the basis of arithmetic, made a considerable use of a mathematical shorthand called algebraic symbolism, and it is very important for the reader, at this early stage, not only clearly to grasp the utility of this symbolism but to make the considerable mental effort required to read it and to perceive its meaning readily. The outstanding advantage of algebra is that it enables us to deal with the questions and topics of mathematics in a general way. Consider, for instance, the test for the divisibility of a number by 9. This is the basis of a large number of arithmetical tricks which appear to indicate that the number 9 has some inherent mystical properties. The divisibility rule, of adding the digits, can be tested for truth by actually trying it on as many numbers as may be desired, and, proceeding in this way, we might say that the more times the test is satisfied the more sure we should be of the truth of the rule. The algebraic proof of the rule which has been given above is, however, quite different from the proof by actual testing. By algebra we are enabled to find the remainder on dividing, not a particular number, but any number by 9. This proof, applicable to all conceivable numbers, is final. Further, and what is almost as important, the algebraic proof shows us clearly that the rule, curious and puzzling at first, has nothing to do with the number 9 of itself, but arises simply because 9 is one less than 10, the basis of our system of reckoning. Thus algebra not only proves a rule in a perfectly general way, but also reveals an important principle underlying it.

These two points are so important that we shall illustrate them by two very simple number games. Consider the old way of discovering a number thought of. The directions are, think of a number, add 3 to it, multiply this sum by 5, add 1 to this product and multiply the sum by 2, and ask the result. If 32 be deducted from this result, the remainder will be a number ending in 0, and, discarding this 0, what is left is the number thought of. Putting the calculation into algebraic symbolism, and calling the unknown number n , adding 3 gives $(n+3)$, multiplying this by 5 gives $5(n+3)$ or $5n+15$, adding 1 gives $5n+16$, multiplying this by 2 gives $2(5n+16)$:

$10n + 32$. This is the result of the calculation for which, if 32 be deducted, a number $10n$ is left. This ends in 0, and when the 0 is discarded n , the number thought of, remains. Set out algebraically this trick, which puzzles some people, is seen to be so simple as almost to be trivial.

Here is a slightly more complicated trick. Take any number of 3 digits of which the first is greater than the third, say 482. Reverse the order of the digits, obtain 284 and subtract, $482 - 284 = 198$; reverse the digits of this last answer and add, $198 + 891 = 1089$. This final answer is always the same whatever number of 3 digits is started with. This can be tested by experiment, but it can be proved easily by algebraic working. Suppose the figures of the number are a , b , and c in order, a being greater than c . The actual number will be $100a + 10b + c$ since a stands for a number of hundreds and b for a number of tens. Reversal of the digits gives a new number $100c + 10b + a$. Let us set the subtraction down in the usual way :

$$\begin{array}{r} 100a + 10b + c \\ 100c + 10b + a. \end{array}$$

Now as a is greater than c , we have to borrow 10 for the units subtraction, and the result is $10 + c - a$, which is a number less than 10. The numbers of tens in the sum are equal in top and bottom lines, but as we borrowed 10 for the units subtraction we have to subtract an extra 10, that is $(b + 1)$ in all, from b tens. We must borrow 100, and we get a remainder of 9 tens. Finally we have to subtract from a hundreds, c plus one borrowed 100; we obtain $(a - c - 1)$ hundreds, so the answer is $100(a - c - 1) + 90 + (10 + c - a)$.

Reversing the digits of this new number gives

$$100(10 + c - a) + 90 + (a - c - 1)$$

and these last two numbers have to be added together. How many units have we? The sum of $10 + c - a$ and $a - c - 1$, or $10 + c - a + a - c - 1$. We add and subtract both a and c , these unknown and indefinite numbers disappear and we are left with 9 units. We see at once that there are 18 tens. The number of hundreds is the same as that of the units, because precisely the same algebraic terms are added together.

The answer is thus 9 hundreds, 18 tens, and 9 units: $900 + 180 + 9 = 1089$. This algebraic working not only proves that the result of the sum is independent of the original number: always 1089; but it shows that this answer depends upon our system of reckoning by tens. For, owing to the borrowing in the subtraction, we obtain in the answer: of units, one less than ten; of tens, twice one less than ten; and of hundreds, one less than ten. The reader who has been interested in this problem might like to try to prove that if we reckoned by any number other than 10, say by r 's, the answer to the sum would be $(r-1) \times (r+1)^2$. The proof is perfectly simple and depends upon putting into algebraic shorthand language the sentence next before the one preceding this.

Review of this Chapter. Our object in this chapter has been to consider the very elementary part of mathematics that is based upon the idea of counting by means of a sequence of cardinal or natural numbers characterised by order, and also to introduce the use of the symbolic language of mathematics whereby, using letter signs, we are enabled to generalise our results and to arrive at assertions, not about particular, assigned, or chosen numbers, but about any numbers. We have seen that, whereas the basis operations of addition, multiplication, and involution are unrestricted in scope, this does not apply to the reverse operations of subtraction, division, and evolution, all of which are inherently restricted so long as we confine our idea of number to the basic counting operation. We have also seen that the symbol 0, first used to denote the emptiness of a class of collections into which a number is divided for the purpose of the decimal system of numeration, behaves as a number in some operations, but not in others. While we have seen that the meaning of some symbols of operation like $-$ and \div is restricted, and that in certain conditions the operations signified are meaningless and inherently absurd, we have encountered one case in which a logical and rational meaning can be assigned to an inherently absurd operation, that of raising a number to a zero power. As our study proceeds we shall see more and more how the development of mathematics consists in finding

logical and rational meanings for those operations which so far have been restricted in scope and absurd outside those restrictions. We shall also see that this leads to ideas about number of a very much farther-reaching character than those based upon the idea of counting.

CHAPTER III

GEOMETRY

We have now to approach mathematics from a standpoint very different from that of the preceding chapter, and to study the elements of geometry, or the branch of mathematics which is concerned with the metrical properties of space. Basically, the properties of space are known only by experience or by the perception of our senses, and as our experience is confined to an extent of space, which is negligible in relation to its boundless character, the study of geometry rests upon certain assumptions about space, which although generally in agreement with experience, cannot be proved in the universal sense.

The earliest systematic development of the science of geometry is found in Euclid's "Elements," a work compiled and arranged with such genius that, up till the end of last century, it was commonly used as an elementary textbook in schools. Although the "Elements" are now rightly regarded, by reason of the strictness and formality of Euclid's exposition, as being unsuitable as a textbook, the work retains its position as one of the greatest that has been produced, and as a lasting monument to the genius of the writer and of the school of classical mathematicians he represented. Although his book has been superseded for educational purposes, Euclid's system is the basis of present-day practical geometrical science.

The science of geometry is commonly said to start with the enunciation of definitions, axioms, and postulates. A postulate is some property of space, unprovable, which has to be taken for granted. Thus Euclid's postulate that a straight line can be produced is really a postulate that space is unlimited in extent. Axioms are common notions about magnitudes, which although unprovable, are self-evident. Definitions are generally descriptive, but sometimes contain an implied postulate. Thus the definition of a square as an equilateral four-sided figure with all angles right, implies a postulate that such a figure can exist.

Straight Lines and Angles and Planes. The conception of a straight line is one which cannot be defined satisfactorily, because a line in the geometrical sense does not exist in the world of sense. A line drawn on paper, however thin, is really a surface, and if the line satisfied the geometrical definition of having no breadth it would cease to exist. So the conception of straightness of a line requires a kind of postulate that there is such a thing.

The conception of surface is almost as difficult to define. It has been said that the surface of, say, a table, is between the table and the air above it—an illustration perhaps as good as can be found. If every line, which can be conceived as lying entirely in a surface, is straight the surface is said to be a plane. It is not sufficient that some lines are straight; an

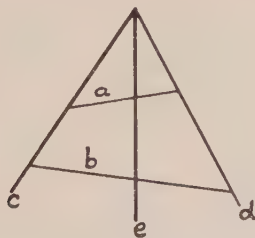


FIG. 1.

indefinite number of straight lines can be drawn on a cylindrical surface, and there is a class of surfaces, apparently curved in every direction, on which an indefinite number of straight lines can be drawn. For a surface to be plane, every line that can be drawn on it must be straight.

If two lines drawn in a plane cross each other or intersect, the part common to each will, if the lines have no thickness, be of no size, and will be merely a position. It is sometimes said that a line is composed of points. So it is, in a way, but the idea leads to some very startling conclusions. Consider the two lines a and b in Fig. 1, terminated by the intersecting lines c and d . Any line such as e , drawn from the intersection of c and d , determines a point on d and a corresponding one on c . Thus for every point on d there is one

on c , so that, although d is longer than c , the numbers of the points of which the two lines are composed must be exactly the same. This paradox may be compared with the one mentioned on p. 25, and it shows that even the most rudimentary ideas about geometry can lead us into deep waters.

When two lines intersect, as in Fig. 2, the inclination between the lines is called the angle they make with each other. Angle is a notion easy to grasp but difficult to define satisfactorily. The opposite angles of the intersection like B and D are evidently equal, but adjacent angles like A and B are generally unequal. If adjacent angles, say, A and B, are equal, as in Fig. 2 (b), the angles are said to be right-angles.

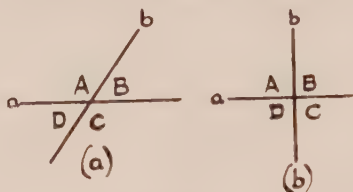


FIG. 2.

What do we mean by angles being equal? It is very difficult to define this exactly, although the meaning of the statement is quite clear. It is a postulate of space that all right angles are equal. The right-angle is therefore a kind of natural measure of angle. Angles less than a right-angle are called acute, those greater, obtuse. It is easy to see that if a number of lines meet another line in one point all the angles at this point on the one side of the last line are together equal to two right-angles.

Parallels. If two straight lines like a and b in Fig. 3 (a) are drawn in a plane at random, then generally these lines will meet and intersect if they are each produced far enough. If a third line, c , is drawn to cross the lines a and b , c is called a transversal, the angles A and B are the inclinations of the lines with the transversal, and these angles are generally unequal. Since the angles A and C make up two right-angles, a and b converge on the side of the transversal on

which the sum of the interior angles B and C is less than two right-angles.

Suppose that, as in Fig. 3 (b), the lines a and b make the same angles with the transversal. Euclid proved that, then, a and b would never meet however far they were produced. Lines which will never meet when produced indefinitely in both directions are said to be parallel, and, it should be realised, that the possibility of parallel straight lines in a plane is a postulate of space. Suppose, however, the two lines a and b converge, by ever so little, so that the interior angles with the transversal are together less than two right-angles; is it certain that they will then meet, if sufficiently produced? Again, supposing that the angles A and B are equal, and the lines

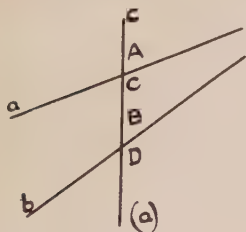


FIG. 3 (a).

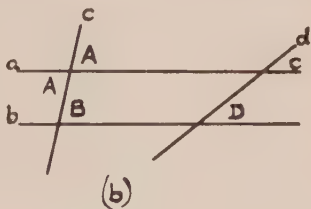


FIG. 3 (b).

parallel, and a second transversal d is drawn; will the angles of a and b with this transversal, C and D , be equal? These two questions cannot be answered definitely, and Euclid had to assume as a postulate that if two lines converged with respect to a transversal, that is, if the interior angles, like B and C , together were less than two right-angles by ever so small amount the lines would meet if produced on the converging side of the transversal. This postulate, when granted answers the second question, for, if the angles C and D were not equal, the lines would converge and meet, which is impossible because they are equally inclined to the transversal c .

Non-Euclidian Geometrics. It is easy to see that Euclid's postulate about parallel lines is equivalent to the statement that through a point one parallel only can be drawn to a straight line. The truth of this seems self-evident to most people, but actually it is not so, and, from the time of Euclid,

many efforts were made to prove it. All these proofs failed, because they begged the question, that is, they were all based on some postulate or other, equivalent to that of Euclid. What is the alternative to Euclid's postulate? Simply this, that from a point A (Fig. 4) it may be possible to draw an indefinite number of straight lines to the left of it, neither of which will ever meet the straight line *a* between the positions given by the line *b*, which is equally inclined with *a* to the transversal *d* and a second line *c* which converges with respect to *c*. Last century two mathematicians, Lobatschewsky and Bolyai, showed separately that, not only is this possibility real, but that a consistent and logical system of geometry can be

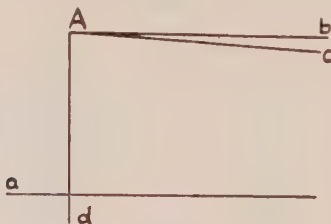


FIG. 4.

built up, by assuming that it is true. A geometry of this kind which rejects Euclid's parallel postulate is called Non-Euclidian. Whether Euclid's, or a Non-Euclidian geometry is true seems hardly capable of proof. Euclid's geometry is by far the simplest, and so far as our sense experience goes, it is certainly true to a high degree of approximation, so that in practical mathematics its strict truth is assumed. It is, however, possible that the Euclidian parallel postulate of space is not exactly true.

We referred to the postulate that underlies the definition of parallel straight lines; that this definition assumes that such lines exist. Riemann showed that this is really an assumption which is unprovable, and he developed a consistent and logical Non-Euclidian geometry in which all lines drawn from a point, like A in Fig. 4, will when produced certainly meet a straight line *b*, that is, in which there are no parallel straight lines.

Non-Euclidian geometries are of no importance in the applications of mathematics to applied science and engineering. Those who wish to obtain a firm grasp of the basis of mathematical theory should, however, realise not only that such geometries are possible, but also, in a general way, something of what they mean.

Triangles. Two intersecting straight lines and a transversal determine a closed figure called a triangle. If the sides or the lines which actually enclose the space within the figure are all equal the triangle is called equilateral; if two sides are equal the triangle is called isosceles. The extent of the surface enclosed by the sides is the area of the triangle. Any three lines will determine a triangle, provided that any two of the lines are together greater than the third. The lengths of the sides of a triangle therefore determine its shape. If two triangles are such that, one being conceived to be placed on the other, they will exactly coincide, they are said to be congruent. If two triangles enclose the same area, the triangles are said to be equivalent. If, with different side lengths, the triangles are of the same shape, that is, corresponding angles are equal, the triangles are said to be similar. The possibility of the existence of triangles, similar but of different size, is a postulate of space which can be shown to be the same as Euclid's parallel postulate, that is, admitting this possibility, Euclid's postulate can be proved.

If two triangles satisfy either of these conditions; corresponding sides are all equal, two corresponding sides and the angle between them are all equal, a side of each and the angles at the ends of these sides are all equal, the two triangles are congruent. This is proved by Euclid but is perceived restrictively. If we conceive two congruent isosceles triangles, the one exactly covering the other, then, if the upper one is turned over it will still exactly cover the other so that the angles at the ends of the equal sides of an isosceles triangle are equal.

Consider the diagram, Fig. 5. ABC is any triangle and a line through A is parallel to BC . AB is a transversal to two parallels so that the angles marked by a cross are equal. AC is also a transversal so that the angles marked by a dot

are equal. Thus the three angles of the triangle are together equal to the three angles under the line at A, that is, to two right-angles. This well-known property of a triangle, it should be noted, rests upon Euclid's parallel postulate, or, its equivalent, that parallels are equally inclined to all trans-

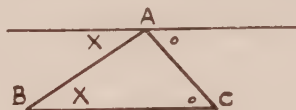


FIG. 5.

versals. In the Non-Euclidian geometry of many non-intersecting lines through a point the angle-sum is less than two right-angles; in the geometry of no parallels, it is greater, and the defect or excess of the angle-sum relative to two right-angles depends upon the size of the triangle.

Areas. A figure enclosed by four lines with all its angles right-angles is called a rectangle. The possibility of the existence of such a figure is a postulate equivalent to Euclid's. The opposite sides of a rectangle are equal, as are those of a four-sided figure called a parallelogram, having opposite sides parallel. If all sides of a rectangle are equal the figure is a square. The size or area of a rectangle is measured by the product of the numbers measuring adjacent sides. The area of a square, the square of the length of its sides, is therefore a square number. Fig. 6 shows that the areas of a parallelogram



FIG. 6.

and a rectangle with a common side and of the same height are equal, for the shaded triangles are plainly congruent. The diagonal or line joining opposite angles of the parallelogram divides the figure into two congruent triangles. Thus, as any triangle is the half of some parallelogram, the area of

a triangle is half that of a rectangle on one side and of the same height.

Circles. A circle is a closed or endless curved line every point of which is equidistant from a point within the space enclosed by the line. This point is the centre. A line drawn from the centre to the line or circumference is a radius, a line through the centre terminated at each end by the circumference, or a double radius, is a diameter, and it divides the circle into two equal and congruent parts. In Fig. 7 BC is

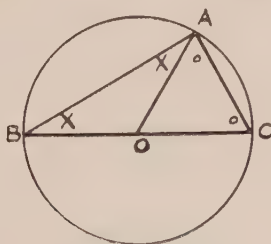


FIG. 7.

the diameter of a circle, and to any point A on the circumference lines are drawn to B, C, and to the centre O. OB, OA, and OC are all equal radii, so that OA divides the triangle BAC into two isosceles triangles. The angles marked with a cross are therefore equal, as are those marked with a dot. The whole angle at A is therefore half the angle-sum of the triangle; this angle is therefore a right-angle. Any angle formed by joining the ends of a diameter to a point on the circumference of a triangle is therefore a right-angle. It can be shown similarly that all angles, like the angle A in Fig. 8,

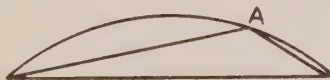


FIG. 8.

formed by a straight line or chord terminated by a part of the circumference of a circle, are equal.

Any circle evidently encloses a definite amount of space

and has a definite area. This area can be divided up into innumerable little triangular areas like that shown in Fig. 9, each triangle having a tiny section of the circumference as one side, and the radius as height. The whole area is then equal to half that of a rectangle with one side equal to the circumference and the other equal to the radius. If the

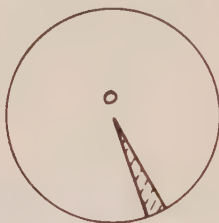


FIG. 9.

measure of the radius is 1 or unity the measure of the area is denoted by a symbol π , which stands for a kind of number which will be discussed more fully in a later chapter.

Ratio, Lines as Magnitudes. The only characteristic of a line, considered by itself is its length. The lengths of two lines drawn at random is generally different and the relation between the lengths of the lines is known as a ratio. Consider the line a in Fig. 10, and the line b , unlimited in length

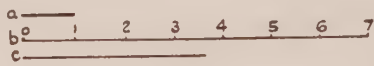


FIG. 10.

towards the right, starting at the point 0, on which are marked a series of points, 1, 2, 3, 4, and so on, such that the distance between any two neighbouring points is equal to the length of a . Suppose a third line, c , is drawn so that one end of it is very near to and opposite the end 0 of b . If the other end of c is exactly opposite a marked point on b , the number of this point indicates how many lines the length a is contained in c . This number is the ratio of the lengths of c and a ; if, for instance, the number is 4, the ratio is 4 to 1, this meaning that 4 times a make up c . Generally, however, the end of c will

not be exactly opposite a marked point on b . The ratio of the lengths of a and c can be expressed by two numbers providing that we can find a length composed of an exact number of length c , which will correspond to the length between 0 and a marked point on b . Thus if 5 times the length of c is exactly equal to the distance between 0 and the point 17 on a , then as 5 c 's are equal to 17 a 's the ratio of the lengths is expressed as 17 to 5. It is evident that, if we can conceive the distance 0 to 17 on b divided into 5 equal parts, each of these parts will be exactly equal in length to c . Hence the ratio of the length can be expressed in the form $\frac{17}{5}$, this meaning that 17 a 's divided by 5 equals c . This kind of division is evidently of a different character from that of repeated subtraction considered in the last chapter, because the length of a line being continuous from start to finish we can conceive that it can be divided into any number of equal parts.

If the length of line a be considered as a unit or equal to 1, then the ratio of the length of c to that of a is the measure of the length of c in terms of the a unit. Practical measurements of length are not made in the manner indicated in the preceding paragraph, but by the use of new units which are exact subdivisions of the principal or primary unit. Thus, in the English system, the inch unit is divided into eight or sixteen equal parts; in the metric system the unit is divided into ten parts. If a were divided decimally, then we should measure c by finding that it contained 3 of the primary units and 4 of the secondary units. If the line c were measured by the English system we should find that no exact number of secondary units, whether the primary unit were divided into 2, 4, 16, 32, and so on, equal parts, would ever exactly equal the piece left over after taking 3 primary units from c . If the ratio of the lengths were 24 to 7, then c could not be measured exactly with either English or decimal secondary units. Provided, however, that the ratio of the lengths of two lines can be represented by two natural numbers, some subdivision of the primary unit will give a secondary unit that will exactly measure one line in terms of the other considered as the primary unit.

There is another way in which we can measure the ratio of two lines, which is interesting in that it does not require a division of the one regarded as a unit. To illustrate this consider lines c and a of Fig. 10, of which the ratio is 17 to 5. Line c contains 3 lengths, a , and a remainder less than a . Let this remainder be applied to a , as a kind of secondary unit, it will be found that a contains two of these first remainders, plus a second remainder less than the first. Let this second remainder be applied to the first, it is found to be contained exactly twice. In this way we measure c by a series of repeated subtractions and arrive at three answers to our divisions, 3, 2, and 2. The ratio can be represented in this way: $\frac{3}{1 + \frac{1}{2 + \frac{1}{2}}}$. The first ratio $\frac{3}{1}$ is the nearest whole number

ratio to the true one and less than it, the next ratio $\frac{1}{2}$ is the nearest whole number ratio of the first remainder to the unit, the last ratio $\frac{1}{2}$ is the exact ratio of the second to the first remainder. This method of obtaining the ratio of the lengths of two lines will be referred to again. We note that whereas it has the apparent advantage, in comparison with the ordinary measuring method, that no subdivision of the unit is required, it does not enable us to solve the converse problem to that of finding a ratio or measure of length, the determination of a line having a given measure or of a length having a given ratio to that of a given unit line.

Generally, it is impossible to express the ratio of the lengths of two lines by two natural numbers in the way we have discussed, because no exact number of times one length or no multiple of it will ever equal any multiple of the second. Otherwise it will be impossible exactly to measure one with the unit by exact subdivision of this unit into any number of equal parts. Ratios of this kind are said to be incommensurable. The ratio of the lengths of two lines is, however, a definite and objective relation between them, whether this ratio is commensurable or incommensurable.

If four lines satisfy the condition that the ratio of the first to the second is equal to that of the third to the fourth, the lines are said to be in proportion. The corresponding sides of similar triangles are proportional, that is, the ratio

of the lengths of similarly situated sets of sides are all equal.

Conceive that squares are drawn on the lines c and a of Fig. 10, which are in the ratio of 17 to 5. If a is divided into 5 parts, 17 of these parts will exactly equal c , the square on c contain 17^2 small squares erected on a fifth part, the square on a will contain 5^2 of these small squares. Hence the ratio of the areas of the squares is 17^2 to 5^2 . The ratio $\frac{17^2}{5^2}$ is the square of the ratio $\frac{17}{5}$. The ratio of the areas of similar figures is the square of the ratio of corresponding sides. Thus the ratio of the areas of circles is the square of the ratio of the diameters, as all circles are similar.

Right-angled Triangles. If one angle of a triangle is a right-angle, the side opposite this angle is called the hypotenuse. As the other two angles together make a right-angle, these angles are said to be complementary. The shape of a right-angled triangle is determined by one of the angles

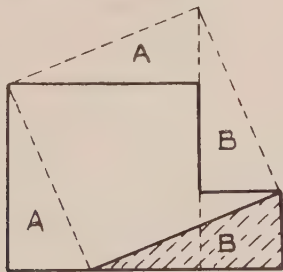


FIG. 11.

other than the right-angle. A remarkable property of right-angled triangles is that the area of the square on the hypotenuse is equal to the sum of the areas of the squares on the other two sides. This is the famous theorem of Pythagoras. Euclid proves this formally, but an informal and convincing proof is given by Dudeney's dissection, illustrated in Fig. 11. The primary right-angled triangle is shaded, and the L shaped full line figure enclosing it is seen to be made up of the squares on the sides forming the right-angle. If the sections of this

figure marked A and B be cut off and transferred to the new positions shown, the dotted square on the hypotenuse of the shaded triangle is obtained. The converse of Pythagoras's theorem is true, that if the sum of the squares on two sides of a triangle is equal to the square on the third side, then the triangle is right-angled.

Consider the numbers 3, 4, and 5. The area of a square with a side of 3 units is represented by $3^2=9$; similarly for the number 4, the square of which is 16. As $16+9=25=5^2$, we see that a triangle having sides of 3, 4, and 5 units respectively will be right-angled. Such a triangle is called a rational right-angled triangle because the lengths of its sides can all be expressed by natural numbers, and a set of numbers like 3, 4, and 5 is called a Pythagorean triplet. A new triplet can, of course, be formed by multiplying the numbers 3, 4, and 5 by the same number; this however is only like magnifying the size of the triangle corresponding. Are there any other right-angled triangles, different in shape from those corresponding to 3, 4, and 5, which are rational? There is an unlimited number; the next simplest is one the sides of which are 5, 12, and 13. $5^2+12^2=169=13^2$. Although it has nothing to do with engineering, it will be of interest briefly to show how Pythagorean triplets can be found, and this will be a good example of the algebraic method of calculation used in the last chapter. If m and n are any natural numbers then $(m^2+n^2)^2$ is easily seen by the argument of p. 14 to be equal to $m^4+2m^2n^2+n^4$. It is also easy to work out similarly that $(m^2-n^2)^2$ is equal to $m^4-2m^2n^2+n^4$. Let us obtain the difference of these two square numbers we have expanded symbolically. Thus difference is

$$m^4+2m^2n^2+n^4-(m^4-2m^2n^2+n^4)=4m^2n^2=(2mn)^2,$$

since, as we saw on p. 13, the double subtraction of $2m^2n^2$ is equal to addition. We therefore obtain a square number as the difference of the two square numbers we started with, hence (m^2+n^2) , $2mn$, and (m^2-n^2) form a Pythagorean triplet whatever values we give to m and n . Let us test the formula and make $m=2$ and $n=1$, $m^2+n^2=2^2+1^2=5$, $2mn=2\times 2=4$, and $m^2-n^2=4-1=3$. This is our first triplet.

Put $m=3$ and $n=2$, and we obtain $9+4=13$, $2 \times 3 \times 2=12$, and $9-4=5$, the 13, 12, 5 set already mentioned. Putting $m=4$, and $n=3$ we obtain 25, 24, and 7. The reader should test these by squaring.

Let us consider the right-angled triangle forming one corner of the large square in Fig. 12. This is half a square, and it is easy to see that in this case Pythagora's theorem is true, for the sum of the squares on the two sides of the triangle, and the square on the hypotenuse are each equal to half the large square. If the sides of the triangle are each one unit, then the square on the hypotenuse has an area of 2 square units, and the ratio of the length of this hypotenuse to the length of the equal sides has a square that is equal to 2. It is easy to show that this latter ratio cannot be

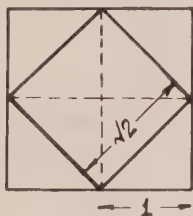


FIG. 12.

expressed in natural numbers, in other words, no square ratio like $\frac{m^2}{n^4}$, where m and n are natural numbers, can represent the ratio 2 to 1. For if it could we should have $2n^2=m^2$. Now m and n cannot both be even or their ratio would be the same as that of two smaller numbers, successively doubled, one of which is odd. As $2n^2$ is even, m^2 is even, and as the square of an odd number must be odd, m must be even. But an even number is the double of some number, $2p$ say, so its square $4p^2$ must be divisible by 4. If $2n^2$ is divisible by 4, n^2 must be divisible by 2, and is even so that n must be even. Thus no natural numbers m and n can be found such that the square of their ratio is exactly 2. The ratio of the hypotenuse to the side of an isosceles right-angled triangle which is half a square cannot therefore be measured exactly by any sub-

division of the unit side. Thus ratio is said to be incommensurable; it is expressed by the symbol $\sqrt{2}$, which is a short way of writing "a ratio whose square is 2." We must note carefully that although it is impossible to express this ratio by natural numbers, it is perfectly definite and objective.

Another important and special right-angled triangle is that shown shaded in Fig. 13. This is exactly half an equilateral triangle, and it is formed by a dividing line which can be proved to bisect the angle through which it passes and to make right-angles with and to bisect the side that it meets. If the small side of this triangle is 1 unit, the slant side or hypotenuse will be 2 units, for all sides of the equilateral

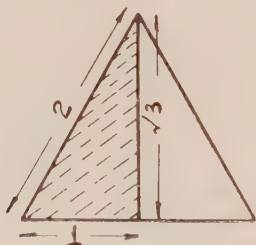


FIG. 13.

triangle are equal. As the square on this slant side has an area of 4 square units, it follows that the square on the vertical side will have an area of 3 units, and the ratio of the length of this side to the length of the horizontal side must have a square equal to 3. This ratio is denoted by $\sqrt{3}$, and like $\sqrt{2}$, it is incommensurable, that is, it cannot be represented as the ratio of two natural numbers.

Trigonometrical Ratios, Angular Measure. Consider the right-angled triangle shown in Fig. 14. The shape of this triangle is determined by the angle A , and all right-angled triangles, of whatever size, having this angle A , have their sides proportional, that is, in the same ratios. The ratio of the side marked a to the hypotenuse c is called the sine of the angle A , this is written $\sin A$ and read "sine A ." The ratio of side b to side c is called the cosine of the angle A , written and read $\cos A$. The ratio of side a to side b is called

the tangent of the angle A, written and read $\tan A$. These ratios are very important, and they should be so memorised that the meanings are recalled instantly. Calling a the

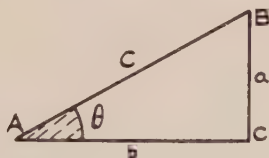


FIG. 14.

height, b the base, and c the hypotenuse of the triangle we can represent these ratios :

$$\sin A = \frac{\text{ht.}}{\text{hyp.}} \quad \cos A = \frac{\text{base}}{\text{hyp.}} \quad \tan A = \frac{\text{ht.}}{\text{base}}.$$

There are other ratios of the sides of the triangle which, although not quite so important, ought to be known. $\frac{c}{a}$ is called the cosecant of A, written and read $\operatorname{cosec} A$, $\frac{b}{c}$ is the secant of A, written and read $\sec A$, and $\frac{b}{a}$ is the cotangent of A, written and read $\cot A$.

As the two small angles of the right-angled triangle make up one right-angle we see, by conceiving the triangle to be turned round so that side a is the base, that the cosine of B is the same as the sine of A. A is the complement of B, and cosine B means the sine of its complement. Similarly, with cosecant and cotangent, as the reader should carefully check.

Suppose now that the hypotenuse c is one length unit, and that b and a are the measures of the sides in terms of these units. a is now $\sin A$, for it is the ratio a to 1; similarly, b is now $\cos A$. But $a^2 + b^2 = 1$, one square unit, the area of c^2 the square on the hypotenuse. Thus

$$\sin^2 A + \cos^2 A = 1$$

$\sin^2 A$, standing for the square of the sine of A. This is an important corollary of Pythagoras which ought to be

thoroughly memorised. Further, as $\frac{a}{b}$ is equal to the $\tan A$ ratio we have

$$\tan A = \frac{\sin A}{\cos A}.$$

We must now inquire a little more closely into the meaning of the statement $\sin A = \frac{a}{c}$. First we see that if this ratio is given in terms of two numbers, say $\frac{7}{11}$, we can at once draw the angle by a fairly obvious geometrical construction; we have simply to make a right-angled triangle with a height 7 of any unit and an hypotenuse of 11 of the same unit. So with $\cos A$ and $\tan A$. Thus $\sin A$, $\cos A$, and $\tan A$ define the angle to which they refer. Further, if one side of a right-angled triangle is given we can, given also $\sin A$, draw the

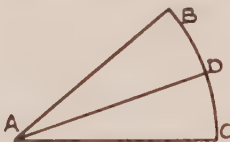


FIG. 15.

triangle. These ratios $\sin A$, $\cos A$, and $\tan A$ are called the trigonometrical ratios because they enable us completely to determine a right-angled triangle (trigon) of which a side is given.

We must however go further than this. It is evident that the ratio $\sin A$ depends upon the shape or size of the angle A ; when A is very small $\sin A$ is a small ratio, when A is nearly a right-angle, the height and hypotenuse of the triangle containing A are nearly equal. $\sin A$ is therefore a kind of ratio associated with the value or size of the angle A . How then can angles be assessed as regards size? In Fig. 15 we see a line AD bisecting an angle A , and the equality of the two parts could conceivably be tested by imagining one part to be put over the other so that it exactly covered it. We might test the equality in another way by drawing a part or arc of a circle to cut the arms of the angle, and by conceiving that the

portions of this are cut off by the two parts of the angle are compared as regards length. Since the circle is a curve of perfect symmetry, we see almost intuitively that the sizes of angles can be considered to be proportional to the lengths of the arcs of a circle they cut off, the centre of this circle being at the point or vertex of the angle.

We can readily see without proof that six equilateral triangles will fit into a circle, as shown in Fig. 16. In the ordinary method of angle measurement the whole circumference of the circle is conceived to be divided into 360 parts, so that the angle of the equilateral triangle contains 60 of these parts or degrees. As all the angles at the centre of a circle make up

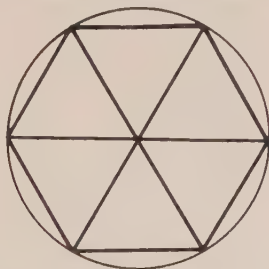


FIG. 16.

four right-angles, one right-angle contains 90 of these degrees. The degree of angle is further divided into 60 minutes (small parts) and each minute can be further subdivided into 60 seconds (second small parts).

Although we may conceive the division of the right-angle into 90 degrees, or, what is the same thing, the whole circumference of a circle into 360 parts, the question arises, can this be done? Euclid showed how, by ordinary geometrical methods the circumference could be divided into 4, 5, 6, and 15 equal parts. Further, as any arc can be bisected, the circle can be divided into numbers of equal parts formed by successive doubling of the numbers 4, 5, 6, and 15. Early last century Gauss proved that an indefinite number of equal subdivisions into parts whose numbers are some primes is possible. Thus, a 17-fold and a 257-fold subdivision is possible by Euclid's geometry of the ruler and compasses. But there is

an infinite number of subdivisions that are impossible, and as these include the division into 360 equal parts the degree measure of angle cannot be realised by simple geometry, that is, although certain angles such as 60, 45, 30, 10 degrees can easily be drawn, most degree values, such as 10, cannot. The equal subdivision of the circumference of a circle is, however, conceivable, and engineers know that it can be achieved mechanically although not geometrically.

There is another method of assessing the size of angles which is of very great importance in mathematical theory. In Fig. 17 an angle A encloses an arc of a circle of unit radius.



FIG. 17.

The size or magnitude of the angle is defined as the ratio of the length of this arc to its radius, and as the measure of the radius is 1, the length of the arc is the measure of the angle. When the length of the arc is 1, the measure of the angle is also 1, and this unit angle is said to be 1 radian. Radian measures of angles are usually denoted by one of the Greek letters θ , ϕ , or ψ , and if we refer to the angle in Fig. 17, not as A, but as θ , this implies that we measure it not in degrees but in radians, that is, as the ratio of an arc to a radius.

There is still another way in which angular magnitude may be regarded, and this is in terms of the whole sectorial area included by the arms of the angle and an arc. By reasoning similar to that used on p. 38 we see that this area is equal to half that of a rectangle having one side equal to the radius and the other equal to the length of the arc. If the radius is 1, then the measure of the length of the arc or radian measure of the angle is equal to twice the sectorial area. Otherwise

radian measure of angle is twice the ratio of a sectorial area determined by the angle to the area of the square on the radius ; if H is the measure of area, and r that of radius, $\theta = \frac{2H}{r^2}$.

As the whole circumference of a circle includes four right-angles, and as the area in the circumference of unit radius is the number π , it follows that 4 right-angles have a circular measure of 2π .

We can now see more clearly, what it is of the greatest importance to grasp firmly, that the trigonometrical ratios are really relations between the lengths of particular straight lines and those of an arc of a circle associated with them. In Fig. 18 we have the right-angled triangle by which $\sin A$,

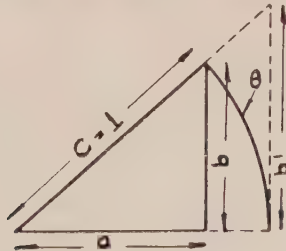


FIG. 18.

etc., were originally defined. The hypotenuse of this triangle, c , is of unit length, and the arc is that of a circle of unit radius. Denoting the angle, represented geometrically as A , by its radian measure or the length of the arc, θ , the statement $\sin \theta = a$ means that when this arc defines a right-angled triangle contained in the radii cutting off the arc, the height of this triangle measures a radius units, $\cos \theta = b$ means that the base of the triangle measures b radius units. It is easy to see from the diagram that if a new right-angled triangle with the radius as base is drawn to enclose the arc the height of this triangle b , is the tangent of θ . Most readers will recognise that this height line touches or is a tangent to the arc, hence the term, tangent of the angle or arc.

Instead of writing $\sin \theta = a$, we can write the equivalent or nearly equivalent statement that θ is an angle or arc whose

sine is a . This inverse statement is abbreviated into the formula, $\text{arc sin } a = \theta$, meaning that θ is an angle whose sine is a . Similarly we may write $\text{arc cos } b = \theta$, and $\text{arc tan } \frac{a}{b} = \theta$.

The symbols arc sin , arc cos , and arc tan all represent angles or arcs defined by their sines, cosines, or tangents, as the case may be. We shall see later, however, that, although the arc completely defines the sine, cosine, and tangent, the converse is not quite true, and that, in the extended idea of angle which we shall study later, there are an indefinite number of angles that can have a stipulated sine, cosine, or tangent.

The reader should make himself thoroughly familiar with this inverse notation, as it is generally found more difficult to read and grasp than the direct notation. What, for instance, does the statement “ $\text{arc sin } a + \text{arc cos } a = 1 \text{ right-angle}$ ” mean, when referred to Fig. 18? Simply that the angle $\text{arc cos } a$, that is, the angle whose cosine is the sine of θ , is complementary to θ . What is the symbolic statement that $\cot \theta$ is

the tangent of the complement of θ ? $\text{Arc tan } \frac{b}{a} + \text{arc cot } \frac{b}{a} = 1 \text{ right-angle}$, or as 4 right-angles make up a radian measure of 2π , this last statement can be put better in the form, $\text{tan } \frac{b}{a} + \text{arc cot } \frac{b}{a} = \frac{1}{2}\pi$. The reader should not leave this paragraph till he has so grasped this inverse notation that he reads it without conscious effort of memory.

An alternative notation for $\text{arc sin } a = \theta$ is $\sin^{-1}a = \theta$. This latter notation, the reason for which will be better understood after the following chapter is studied, is generally considered undesirable, although used considerably in this country. In time it may become obsolete.

Vectors. Hitherto we have considered straight lines merely in reference to their lengths and to the ratio these lengths bear to a unit length. The length of a line can, it is evident, stand for a movement from one end of the line to another, or for an objective distance. Thus a straight line 3 inches long drawn anywhere on a 1-inch Ordnance map will represent a distance of 3 miles, and, starting at the physical position at one end

of the line, it will show the new position arrived at by a 3-mile movement as the crow flies. The line on the map will, however, show more than this, for as vertical edges of the map correspond to north and south and horizontal edges to east and west directions, the inclination of the line to these edges or boundaries will indicate a direction. If the line is inclined at 45 degrees to the lower boundary of the map it will indicate a direction of movement which will be NE or SW accordingly as we start from the lower or upper termination of the line. We can remove this ambiguity as regards the particular direction to be indicated by drawing an arrow-head at one end of the line. If this arrow-head is at the upper end the line then stands definitely for a 3-mile movement in a north-easterly direction.

A straight line which is used in this way to indicate completely a specified kind of movement, or, as it is sometimes called, a directed step, is called a vector. A vector has three basic characteristics, first a length, representing actually or to some scale, an extent of a movement or step, second a direction relative to some other fixed line (an edge of the map in our illustration) which indicates angular inclinations and, thirdly, an arrow-head which denotes and finally decides direction or sense.

In Fig. 19 we have shown what might be a line *a* drawn on a map to indicate a directed step in a north-easterly direction, of an amount represented by the length of the line. Although a straight north-easterly movement might be made, as indicated, by a bird, it would rarely be possible, excepting in quite open country, on the surface of the earth. If we wished to proceed on foot from the position indicated by the start of the line or vector to its end we might, even using field-paths, have to take a devious route indicated by *b*, and our actual movements would be, first one nearly north, and then one nearly east. Provided, however, our movements were indicated by the two straight directed steps of the *b* route we should arrive at the place indicated by the termination of the *A* vector. The vector *A* is therefore in a sense the sum of the vectors comprising the *b* route, for the result of the two directed steps of the *b* route gives a resultant or net displace-

ment the same as is indicated by the original A vector. If our journey had to be made by car it might be necessary to follow a still more devious route indicated by c . In this case we should almost have to retrace our path in the last stage of the journey, but the A vector would still be in a sense the sum of the three vectors composing the c route, because the net movement as the result of proceeding along both routes, the

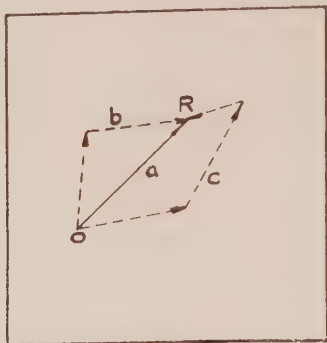


FIG. 19.

one a direct and the other c devious, would be exactly the same.

It is evident that in saying that vector A is the sum of the three vectors in c route, or is the result of adding these vectors, we are using the word addition in a sense very different from that in which it was used in Chapter II. We shall elucidate this extended meaning of the word addition in the chapter following.

Suppose we wish to specify definitely how a vector of a directed step is to be drawn on a map from a stipulated point, such as O in Fig. 20. We can do this in two ways. The first and most obvious is to specify the length, r , of the vector, and the angle θ it makes with one of the cardinal directions on the map. Thus at Fig. 20 (a) the vector would be specified as of length r and inclined at an angle θ , northward to the easterly direction OE . Similarly, at Fig. 20 (b) the direction specification would be an angle southward to the westerly direction OW . There is another way in which this specification could be

made. The vector in Fig. 20 (a) is the resultant or sum of two vectors, the one eastward of length x and the other northward of length y . The r vector therefore would be completely defined as equivalent to the sum or resultant of an x easterly and a y northerly vector. Similarly the vector in Fig. 20 (b) is completely defined by an x westerly and a y southerly vector. These two methods of specification are of great importance in mathematics; the first by length and angle is called the Polar method, the second by lengths in the cardinal

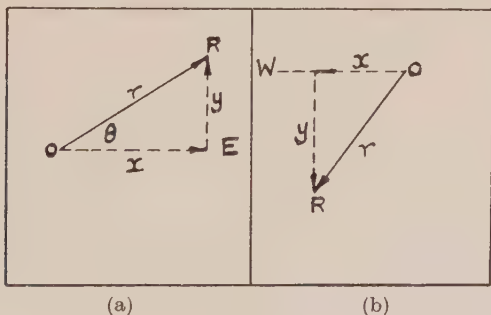


FIG. 20.

directions is called the Cartesian method. The importance of these methods will be seen hereafter.

The characteristic properties of a vector, x and y in the Cartesian specification and r and θ in the polar specification, are evidently not independent. We see at once by Pythagoras that as r , x , and y form a right-angled triangle, $r^2 = x^2 + y^2$. Further, the ratio $\frac{x}{r}$ is the cosine of the angle θ , and the ratio $\frac{y}{r}$ is equal to $\sin \theta$.

We shall see in the following chapter how, by extension of the meaning and scope of the idea of number, these methods of the specification of a vector are simplified, and how a vector itself becomes a new kind of number.

CHAPTER IV

NUMBER

Counting and Measurement. In the two preceding chapters we have considered the ideas underlying the two basic practical applications of mathematics, counting and measurement. Counting a collection of objects consists essentially in placing them in one-to-one correspondence with the series of natural numbers, beginning at the first. The cardinal number of the collection is the last natural number associated with an object of the collection. Measurement of the length of a straight line arises from the idea of ratio, the quantitative relation existing between the length of this line and that of another line chosen as a unit. The primary idea of a number-statement of ratio depends on successive multiplications of the line and the unit till two multiples are found that are equal. But this is the same as the subdivision of the unit line into such a number of equal parts that some exact number of these parts will make up the line to be measured. The result of measuring a line in this way is expressed by the symbol $\frac{a}{b}$, a and b being natural numbers, of which b denotes the number of equal parts into which the unit is divided and a the number of these parts that has to be taken exactly to equal the measured line.

The measurement of a line in this way involves counting and is of very wide application. * We have already seen that angular measurement depends upon the measure of the length of an arc, although we have not considered the question how the length of a curved line can be compared with that of a unit. Surface and volume measures depend upon length measures. The accurate determination of weight in a laboratory depends partly on counting standard weights, and ultimately on the measurement of the distance of a rider on the balance arm from the point of suspension of this arm. Temperature measurement depends upon the determination of the length of a mercury column, and most electrical measurements ultimately depend upon either length or angle measurements.

Counting and measurement are, however, essentially different in character. The operation of counting is used to determine what is called discrete magnitude, magnitude that varies only in definite steps. If there are five people in a room, the number of persons can be increased at the least by exactly one; a line of length 5 inches can, however, be conceived gradually to grow till it attains the length of 6 inches, and it can, unlike the collection of persons, have any magnitude between 5- and 6-inch units. Counting is exact, but measurement of continuous magnitude can, in practice, be only approximate. The reckoning of money is a counting operation, so that an error of 1*d.* in a book balance is as important to the accountant as an error of £100.

Fractions. The statement of a ratio by two numbers, like $\frac{5}{7}$, is called a fraction, because it depends upon the idea of the breaking up of a unit into an equal number of parts, 7. This number is called the denominator of the fraction—it indicates its primary character; the other number, 5, indicating the number of equal sevenths to be taken, is called the numerator. Fractions of the same kind, or having the same denominator, can be added in the ordinary way; $\frac{5}{7} + \frac{4}{7}$ means the sum of a collection of 5 + 4 sevenths, and gives a result 9 sevenths, or $\frac{9}{7}$. Similarly with subtraction. The denominator of a fraction can be changed by further subdivision of each part represented by this denominator. Thus, as $\frac{1}{7}$ is equal to $\frac{4}{28}$, $\frac{5}{7}$ will be equal to $\frac{20}{28}$. Similarly, as one-quarter, or $\frac{1}{4}$, is plainly $\frac{2}{8}$, $\frac{3}{4}$ is $\frac{6}{8}$, so that we have $\frac{5}{7} + \frac{3}{4} = \frac{20}{28} + \frac{21}{28}$, and, as 20 + 21 give a total of 41, the result of the addition is $\frac{41}{28}$. We can express this process of adding fractions in a general way by algebraic symbolism; $\frac{a}{b} + \frac{c}{d} = \frac{ad + bc}{bd}$. In the numerator a is multiplied by d and c by b to make the two fractions of the same kind, that is, so that they each have the same denominator, bd .

The multiplication of a fraction by a natural number is quite a straightforward operation. $\frac{4}{7}$ multiplied by 3 plainly means 4 × 3 sevenths, or $\frac{12}{7}$. This operation ought to be commutative if the multiplication of fractions is to be like that of numbers, but what meaning are we to attach to the

process of multiplying 3 by $\frac{4}{7}$? Although this process is evidently different from repeated addition, a fairly obvious meaning can be given to it; divide 3 into 7 equal parts, each of which will evidently be equal to $\frac{3}{7}$ of the parts obtained by subdividing the unit into 7, and take 4 of these parts. The result is $\frac{3}{7} \times 4$ or $\frac{12}{7}$, the same as was obtained by multiplying $\frac{4}{7}$ by 3. Thus the operation of multiplying a fraction by a natural number is commutative. The operation of multiplying a fraction by a fraction can be similarly interpreted: $\frac{4}{7}$ multiplied by $\frac{3}{5}$ means subdivide $\frac{4}{7}$ into 5 equal parts, making the ultimate division of the unit into 7×5 parts, and giving $\frac{4}{7 \times 5}$; and then take 3 of these parts, giving $\frac{4 \times 3}{7 \times 5}$. This operation can easily be seen to be commutative, and we can say, generally, that $\frac{a}{b} \times \frac{c}{d} = \frac{ac}{bd}$.

The division by fractions is somewhat more subtle. Consider first the division of a natural or whole number by a fraction consisting of 1 part. $1 \div \frac{1}{7}$ means how many times can $\frac{1}{7}$ be taken or subtracted from 1. Now that we can conceive 1 to be subdivided, the answer is plain; 7 times. How many times can $\frac{1}{7}$ be taken from 2? plainly 14 times. What is the result of dividing 2 by $\frac{3}{7}$? As the divisor is 3 times as large the answer ought to be one-third of that obtained by division by $\frac{1}{7}$, that is $\frac{14}{3}$. Arguing in this way it appears that, symbolically $a \div \frac{b}{c} = \frac{ac}{b}$. Having established this, the division of a fraction by a fraction can be interpreted, for, if dividing a by $\frac{b}{c}$ gives $\frac{ac}{b}$, dividing a smaller number $\frac{a}{d}$ by the same fraction $\frac{b}{c}$ will give the original answer divided by d , that is, $\frac{a}{d} \times \frac{c}{b}$ or $\frac{ac}{bd}$.

The operation of division, now that we can use fractional numbers, is plainly less restricted than it was in Chapter II. Thus, originally, in dividing 21 by 5, we had to give the answer as 4 and 1 over, but we can, using fractions either express this division directly as a fraction $\frac{21}{5}$, or we can split up 21 into two parts, 20 and 1, divide each by 5 and give the answer as $4 + \frac{1}{5}$ which is contracted into $4\frac{1}{5}$.

There is plainly no limit to the number of fractions between successive natural numbers because between any fractions, however close together, we can always insert another by adding them together and dividing the result by 2. We can return to a geometrical consideration of fractional numbers. Conceive a limited straight line. If this is divided into 2 parts the point of subdivision will be part of the line. Subdivision into 3 will give 2 more points; into 4, 2 more; into 5 will give 4 more, and so on. This kind of subdivision can go on indefinitely, but however far we go the points we obtain will never make up the whole line; they will remain detached and separate. But as the distances of the points so obtained from the end of the line represent fractional numbers, we feel that the distances of the other points in the gaps or empty spaces, between these fractional number points, ought also to represent some kind of numbers. We have, as a matter of fact, already encountered one of these anomalous numbers which we have indicated by the symbol $\sqrt{2}$. The diagonal of a square of 1 inch side is between 1 inch and 2 inches in length. If we conceive a line exactly 2 inches long on which is a point distant from one end by exactly the length of the diagonal of a 1-inch square, then, into however many equal parts this 2-inch line may be divided, no point of equal subdivision can ever coincide with the point representing the length of the diagonal, for if it could it would be possible to represent the measure of the diagonal as a fraction or exact ratio, which we showed on p. 43 to be impossible. We shall consider the anomalous numbers suggested by these incommensurable lengths in greater detail hereafter.

Decimal Fractions and Percentages. Fractions so far considered are called vulgar (common) fractions. A special class of fractions in common use are those in which the denominators are powers of ten, and in which these denominators are not shown explicitly, but are indicated by the position of the numerator relative to a mark or decimal point, in the same way as the denomination of a figure—units, tens, hundreds, and so forth, is indicated. As in the number 23 each unit of the set of 3 stands for something, one-tenth of what is meant by each of the set of 2, so in the number 23·4, the 4 stands

for a number of subdivisions of the ones indicated by the 3, each of these subdivisions being one-tenth of 1. So 23·4 denotes $23 + \frac{4}{10}$, and 23·45 would denote $23 + \frac{4}{10} + \frac{5}{100}$, or as $\frac{4}{10}$ is the same as $\frac{40}{100}$, $23\frac{45}{100}$. The part of a number of this kind which is separated by a mark or decimal point is called a decimal fraction, and it is equal to a fraction of which this part of the number is numerator and of which the denominator is a 1 followed by as many zeros as the number of figures in the numerator. When a number consists of a decimal fraction only, a zero is often placed in the blank units position to localise the decimal point. Thus ·45 is often written 0·45. As with a natural number, some of the places standing for certain fractions may be vacant; 0·405 means that there are no $\frac{1}{100}$ th fractions in the number. 0·045 means that there are no tenths. Decimal fractions can be added, subtracted, multiplied, and divided by the ordinary rules of arithmetic provided certain precautions are taken which are explained in the textbooks. Multiplication and division of numbers composed wholly or in part of decimal fractions is made straightforward by expressing these fractions in the common or vulgar form. Thus $12 \times 0\cdot0007$ is the same as $\frac{84}{10000}$, and as there are four zeros in the denominator there must be four figures in the decimal fraction, hence the answer is 0·0084. Similarly 0·4 divided by 0·002 is the equal to $\frac{4}{10} \div \frac{2}{1000}$ or $\frac{4}{2} \times \frac{1000}{10}$ or to 200.

Some vulgar fractions can be converted into decimal fractions. Thus $\frac{1}{4}$ is found by dividing 4 into 100 to be $25\frac{1}{100}$ ths or 0·25, $\frac{1}{25}$ can contain no $\frac{1}{10}$ th, but is by division found to be $4\frac{1}{100}$ ths, and hence is written decimally as 0·04. The technique of converting vulgar into decimal fractions is given in textbooks on arithmetic.

Most vulgar fractions cannot be converted exactly into decimal equivalents. Thus $\frac{1}{3}$ is not an exact number of tenths, hundredths, thousandths, or of any decimal fraction however large the denominator. The ordinary process of converting $\frac{1}{3}$ gives an answer 0·333 . . . which is an endless succession of 3's, endless because however far we carry the division by reckoning in smaller and smaller decimal fractions we always have a remainder of 1. Decimal fractions of this kind are

called circulating and are indicated by dots placed over the circulating or recurring figures. $0.\dot{3}$ means that the 3 is repeated endlessly. $0.1\dot{6}$ means that after the 1 the 6 recurs. Circulating decimals are of no importance in practical calculations, and there is no need to give here the rules for manipulating them, although later on we shall see a little more of their significance. No vulgar fraction can be converted into an exact or terminating decimal fraction, unless its denominator is made up of products of 5's and 2's; if any other factor occurs the decimal will circulate. It is a useful mental exercise to try to see why this is so.

Percentages are a particular kind of decimal fractions in which the basic reckoning is in hundredths. Thus 0.04 being 4 hundredths can be expressed as 4 per cent. So also 0.025 is 2.5 or $2\frac{1}{2}$ per cent. A decimal fraction can evidently be converted into a percentage by moving the decimal point two places. The percentage notation is useful because it gives a more definite and less easily mistaken way of writing small fractions; 3 per cent. is perhaps less easily mistaken than .03. On the other hand, the indiscriminate use of the percentage notation gives occasion for that avoidance of simplicity which is a defect of much speaking and writing. As it is usual to say commence instead of begin, check up on instead of check, unilateral instead of one-sided, so it is becoming increasingly common to say thirty-three and a third per cent. instead of one-third, and for engineers to refer to 100 per cent. of full load instead of, simply, to full load.

Approximation. Reference has already been made to the fact that measurement can never be more than approximate. The great advantage of the decimal notation for fractions is that in the numerical statement of a measure it can be made to indicate, not only the actual measure, but also the degree of exactness or approximation with which it has been made. Thus the statements of a length in inches as 2, 2.0, 2.00, and 2.000 are numerically identical, the 0's to the right of the decimal point and the point itself are all unnecessary and redundant if the length is exactly 2 inches. Actually, in applied mathematics, these numbers have all different meanings. To express a length in inches as 2 means, or ought to

mean, that it has been estimated by eye, that is, that it is less than $2\frac{1}{2}$ and greater than $1\frac{1}{2}$ inches ; to state the length as 2.0 means that, in the measurement, tenths of an inch have been taken into account and that the length is less than 2.05 and greater than 1.95, the sort of measurement that would be made by an ordinary foot-rule ; 2.00 would imply that the measurement has been made by a vernier calliper in which $\frac{1}{100}$'s of an inch can be accurately estimated ; 2.000 implies a method of measurement by which it is certain that the length does not differ from exactly 2 by more than half a mil or thousandth of an inch.

The amount by which a magnitude can, due to the imperfections of the measuring process, differ from the recorded number may be called the possible error, and it is often useful to express this as a fraction. If a measure is recorded as a , and the possible error is α , then the fractional possible error $\frac{\alpha}{a}$ is a kind of measure of precision of the measurement. Thus a length stated as 2.11 inches implies that the length lies between the limits 2.115 and 2.105, and the fractional possible error is $\frac{0.005}{2.11}$, and this error is deemed to be of negligible importance. The uncertainty of the measurement can be expressed symbolically as $a\left(1 \pm \frac{\alpha}{a}\right)$, the symbol \pm denoting that the possible error may be in excess or in defect.

In the preceding paragraph we have used to denote a possible error of a so small that it can be neglected by the corresponding Greek letter α . There is another and a better notation for a small quantity of this kind with which the reader must be thoroughly acquainted. According to this notation the quantity α would be written δa , the δ here being not the symbol for a number but for a mark, attached to the a , showing that it is some small uncertainty in the value of a , or some small difference (hence the Greek letter δ) between the stated and the actual values of the magnitude referred to. There is no doubt that the possible ambiguity of the whole symbol δa and the possibility of this meaning the multiplication of two magnitudes $\delta \times a$ is one of the defects of

mathematical symbolism, but actually the symbol δ is rarely used in mathematics in any other way than as a mark of a small difference, and once the notation has become familiar there will hardly ever be any actual ambiguity in the meaning of the symbol. Thus, using this symbol, we may say that the actual magnitude of a quantity, measured as a , may vary between the two limits contained in the formula $a\left(1 \pm \frac{\delta a}{a}\right)$,

this being a condensation of $a\left(1 + \frac{\delta a}{a}\right)$ and $a\left(1 - \frac{\delta a}{a}\right)$. The fractional possible error $\frac{\delta a}{a}$ is often expressed as a ratio, thus

a length noted as 2.8 inches, lying possible between 2.85 and 2.75, might be said to have a possible error of 5 parts in 280.

Let us now consider the multiplication of two approximate numbers like $a\left(1 \pm \frac{\delta a}{a}\right)$ and $b\left(1 \pm \frac{\delta b}{b}\right)$. Note carefully that δb stands for the possible error in the b measurement. The true product can evidently lie between the product of the two greatest and the two smallest values of the factors. This first product is $a\left(1 + \frac{\delta a}{a}\right) \times b\left(1 + \frac{\delta b}{b}\right)$, or $ab\left(1 + \frac{\delta a}{a}\right)\left(1 + \frac{\delta b}{b}\right)$. If we work out the expansion of the product of these two quantities in brackets in the way explained on p. 14 we find the answer to this sum to be $ab\left(1 + \frac{\delta a}{a} + \frac{\delta b}{b} + \frac{\delta a \times \delta b}{a \times b}\right)$.

Now $\frac{\delta a}{a}$ and $\frac{\delta b}{b}$ are both fractions small enough to be neglected, so that their product will be a much smaller fraction. Thus if the fractional errors are each 0.01 their product will be a much smaller fraction, 0.0001. We can therefore neglect this product and say, that to the degree of approximation we are considering, the greatest value of the product of two magnitudes measured as a and b will be $ab\left(1 + \frac{\delta a}{a} + \frac{\delta b}{b}\right)$. The possible fractional error of the product is therefore the sum of the possible errors of the factors. We obtain a similar result if

we multiply the least possible true values of the measured magnitudes.

Let us now see what this means. Suppose we multiply in the usual way the two approximate numbers 1.04, and 1.21, the answer is 1.2584. The possible fractional errors in the factors are 5 parts in 1040 and 5 parts in 1210, each roughly 5 parts in 1000. The possible error in the product is therefore 10 parts in 1000, or 1 part in 100. This means that the product of the true values may be as much as 1 part in 100 in excess or in defect of the answer to the sum, 1.2584; in other words, it may be as much as $1.2584 + 0.0125$, or 1.2709. We see then that the two right-hand figures 8 and 4 in the calculated product are quite meaningless, and that even the 5 is uncertain. If we work out symbolically the division of one approximate number by another we shall obtain a similar result, for the answer to the sum $a \div \left\{ b \left(1 + \frac{\delta b}{b} \right) \right\}$ multiplied by $b \left(1 + \frac{\delta b}{b} \right)$ must equal a , so that this answer must be $\frac{a}{b} \left(1 - \frac{\delta b}{b} \right)$; for multiplying this by $b \left(1 + \frac{\delta b}{b} \right)$, we obtain $a \left\{ 1 + \frac{\delta b}{b} - \frac{\delta b}{b} - \left(\frac{\delta b}{b} \right)^2 \right\}$, which, as the square of the fraction is so small that it can be neglected, is equal to a .

Continued Fractions. We have already seen on p. 40 that a method of measurement of length alternative to that of using a subdivided unit is to find the numbers of times that first, second, and successive remainders are contained in the unit and in remainders of the next higher order. Thus, if the unit is contained in the original length twice, the remainder in the unit three times, the second remainder in the first four times, and the third remainder in the second twice, we can express the measure of the length as $2 + \frac{1}{3 + \frac{1}{4 + \frac{1}{2}}}$. A continued fraction gives successive approximations or convergents to the actual value it represents, and these convergents are alternately lower and higher than the actual value. Thus, considering the fraction given, the first approximation or convergent is 2, obviously too small; the next is $2\frac{1}{3}$ obtained by

regarding the first two members of the fraction only, and this is too large because the denominator of the fraction taken as 3 is actually larger than 3; in the third convergent we take the denominator of the fraction as $3\frac{1}{4}$ or $\frac{13}{4}$, so that the convergent is $2\frac{4}{13}$, a little too small, because the denominator of $\frac{13}{4}$ is too small; the true denominator of the fraction is not $3\frac{1}{4}$ but 3 plus 1 divided by $4\frac{1}{2}$, or $3\frac{2}{9}$; this is equal to $\frac{29}{9}$ so that the final value of the whole fraction is $2\frac{9}{29}$. The serious reader should carefully follow this argument.

Continued fractions are of considerable importance in pure mathematics, but do not concern the engineer very much. Reference should however be made to one use of continued fractions in practical calculation; this is the expression of a decimal fraction in a simpler form as a vulgar fraction. The approximate value of π , the length of the circumference of a circle of unit radius can be shown by calculation to be 3.14159, or $3\frac{14159}{100000}$. If 14159 is divided into 100,000, it is contained 7 times with a remainder 887; 887 is contained 15 times in 14159 with a remainder 844, and 844 is contained almost exactly once in 887. 3.14159 is therefore very nearly equal

to $3\frac{1}{7+\frac{1}{16}}$. The denominator of the fraction following

3 is thus $7\frac{1}{16}$ or $\frac{113}{16}$, so that the complete number is $3\frac{113}{16}$ or $\frac{355}{113}$. The reader should check by ordinary arithmetic that this fraction is, to five places of decimals, equal to 3.14159.

Irrational Numbers. Any number which can be expressed as a combination of operations of adding, subtracting, multiplying and dividing natural numbers is called a rational number, and we have seen that, between any two natural numbers like 1 and 2, there is an unlimited number of rational numbers. We have also seen that in a line 2 inches long there is a point representing a length, the square on which has an area exactly double of that having a 1 inch side, and we have represented this length symbolically as a new kind of number in the form $\sqrt{2}$. Further, we have seen that the height of an equilateral triangle described on the 2-inch line can geometrically be said to have a length of $\sqrt{3}$ inches. The new number $\sqrt{2}$ is not a rational number, for it cannot be expressed as a fraction if we consider that such a number can be multiplied by

itself or squared, then we can say the $\sqrt{2}$ is such a number that has a square equal to 2. There is evidently an assumption or postulate in this statement that a number, which cannot be represented by the division of one natural number by another, can be multiplied by itself in the arithmetical sense. If we make this assumption we can regard the number represented by $\sqrt{2}$ as a new kind of number, of which we know that its square is exactly 2. $\sqrt{2}$ is called an irrational number because it cannot be expressed in a rational form. Admitting our multiplication postulate we can say that $\sqrt[3]{2}$ is another irrational number characterised by the property that $\sqrt[3]{2} \times \sqrt[3]{2} \times \sqrt[3]{2} = 2$. The irrational numbers comprise all square roots of numbers that are not exact squares, all cube roots of numbers that are not exact cubes, and so on, together with other numbers resulting from the addition, subtraction, multiplication, and division of these elementary irrationals. It follows therefore that the irrational numbers between, say 1 and 2, are much more numerous than the rational numbers, for they include all roots of all orders of every rational number in the interval 1 to 2.

The irrational numbers are a class, continuously mingled with the rationals, but in a way standing apart from them. The sum of an irrational, say \sqrt{a} plus a rational b , cannot equal a second rational number c , for if this were the case \sqrt{a} would be equal to $c - b$ and would be rational. Irrational numbers multiplied together usually give a third irrational number. We can assign a meaning to this idea of multiplication in the following way. Consider $\sqrt{a} \times \sqrt{b}$; if this represents a new number let this new number be squared, the result is $\sqrt{a} \times \sqrt{b} \times \sqrt{a} \times \sqrt{b}$ or $\sqrt{a} \times \sqrt{a} \times \sqrt{b} \times \sqrt{b}$ or $a \times b$. The square of the product is ab so that $\sqrt{a} \times \sqrt{b} = \sqrt{ab}$. This should not be regarded as too obvious. Sometimes the product of two different irrationals seems to give a rational answer, thus $\sqrt{8} \times \sqrt{2} = \sqrt{16} = 4$. If, however, we consider that $\sqrt{8} = \sqrt{(4 \times 2)} = \sqrt{4} \times \sqrt{2} = 2 \times \sqrt{2}$, we see that the multiplication of irrationals really only involves the squaring of $\sqrt{2}$, which, of course, gives a rational result. We may note here that a symbol like \sqrt{a} is conventionally meant to represent an irrational number; this qualification is necessary

because, if a were considered as any number whatever, it might have such a value as 4, in which case \sqrt{a} would be rational.

A large and important class of irrational numbers in practical mathematics are those involving square roots, and these numbers are sometimes called quadratic irrationals. Another name is quadratic surds, the name "surd" being used to indicate an old view that, as incommensurable ratios could not be added or multiplied by the processes of common arithmetic, it was absurd to consider them as representing numbers. A fraction having for a denominator a quadratic surd can, by a simple artifice, always be transformed to an equivalent form, more convenient for calculation. $\frac{1}{\sqrt{a}}$ is

unchanged if numerator and denominator are multiplied by the same quantity \sqrt{a} , thus $\frac{1}{\sqrt{a}} = \frac{\sqrt{a}}{\sqrt{a} \times \sqrt{a}} = \frac{\sqrt{a}}{a}$, and $\frac{1}{\sqrt{2}}$ can be expressed as $\frac{1}{2}\sqrt{2}$. This process is called rationalising the denominator. Again, consider the fraction $\frac{1}{a + \sqrt{b}}$. If

we multiply $(a + \sqrt{b})$ by $(a - \sqrt{b})$, this is multiplying the sum and difference of two numbers, the result of which, as we saw on p. 14, gives the difference of their squares. Hence $\frac{1}{a + \sqrt{b}} = \frac{a - \sqrt{b}}{(a + \sqrt{b})(a - \sqrt{b})} = \frac{a - \sqrt{b}}{a^2 - b}$, a fraction with a rational denominator. In this calculation $(a - \sqrt{b})$ is called a rationalising factor.

Approximate Values of Irrational Numbers. As the length of the diagonal of the side of a unit square can be determined approximately by actual measurement, so an approximation to an irrational number can be obtained in the form of a rational number, which has nearly the same characteristic property as the irrational. The books on arithmetic give a method of calculating approximate values, for instance, of $\sqrt{2}$; the answer to a sum of this kind is 1.41421, to five decimal places. What do we mean when we say that this rational number 1.41421 is nearly equal to $\sqrt{2}$? Simply and only this, that if we square 1.41421 by the rules of arithmetic we obtain an answer that is very nearly equal to 2. We can continue

the arithmetical calculation, and obtain as many decimal places in our answer as we please, and the more we obtain the closer we may consider the approximation, because the square of our answer will become nearer and nearer to 2. The process of extracting a square root as it is called, or of obtaining a rational approximation to an irrational, if continued indefinitely will never cease, and the figures of the decimal fraction of the answer can never recur, for, if they did, the unlimited line of figures would represent some rational fraction. The rule for obtaining an approximation to an irrational square root is easily learned but easily forgotten, but it is rarely required by the engineer. Rules have been devised for obtaining irrational cube and fifth roots approximately, but these are more difficult to manipulate. Theoretically it should be possible to devise a rule for the finding of an approximate value of any irrational root, but for roots higher than 5 which are expressible by prime numbers the calculation would be insuperably difficult. Such approximate roots, if required, are calculated indirectly.

Although a quadratic surd or an irrational involving a square root cannot be expressed exactly as a recurring or circulating endless decimal fraction, it is of interest to note that it can be expressed as a recurring continued fraction. We can easily show how $\sqrt{2}$ can be put into the form of a recurring continued fraction for $\sqrt{2} = 1 + (\sqrt{2} - 1)$, and $(\sqrt{2} - 1)$ can be turned into a fraction by the reverse process of rationalising a denominator for it is equal to $\frac{2-1}{\sqrt{2}+1}$ or $\frac{1}{\sqrt{2}+1}$. The denominator of this fraction can be put equal to $2 + (\sqrt{2} - 1)$ so that $\sqrt{2}$ is equal to $1 + \frac{1}{2 + (\sqrt{2} - 1)}$. We can now treat $(\sqrt{2} - 1)$ in the same way as we did before, and obtain $1 + \frac{1}{2 + \frac{1}{2 + (\sqrt{2} - 1)}}$, and this can be continued endlessly. We therefore see that $\sqrt{2}$ can be considered to be equivalent to the continued fraction $1 + \frac{1}{2 + \frac{1}{2 + \dots}}$, recurring indefinitely, in a way similar to that in which $\frac{1}{3}$ can be considered to be equivalent to the endless decimal fraction 0.333 . . .

It is often useful to know very roughly what the root of a large number is, say 1,345,211. A question of this kind can easily be answered by expressing the number as a multiple of some power of 10, the required root of which can easily be obtained. 1,345,211 is roughly 1.3 millions or 1.3×10^6 and the square root of this is the product of the square root of its factors. As the square root of 1.3 is a little over 1, and the square root of 10^6 is obtained by dividing the index by 2, the required rough square root is a little over 1000. If we required the rough fifth root we should have to express the number as a multiple of some power of 10 having an index divisible by 5. This would be as $13 \times 100,000$ roughly or 13×10^5 . The fifth root of 13 is less than 2, and the fifth root of 10^5 is simply 10, so that the required fifth root lies between 10 and 20.

Fractional Indices. The recognition of irrational numbers evidently widens the scope of evolution, which was discussed in a rudimentary way on p. 25. As long as the idea of number is confined to the succession of cardinal numbers $\sqrt[n]{a}$ has no meaning excepting when it is a perfect square. We have seen that in the restricted sense of evolution, the n th root of a^m can be found only when m is exactly divisible by n in which case it is $a^{\frac{m}{n}}$. Let us remove this restriction and see if there is any meaning that can be assigned to an index which is a fraction inherently, and not a whole number. What, for instance, can be the meaning of $a^{\frac{1}{2}}$. By the rule of addition of indices $a^{\frac{1}{2}} \times a^{\frac{1}{2}}$ ought to equal $a^{\frac{1}{2} + \frac{1}{2}} = a^1 = a$, so that it is reasonable to consider that $a^{\frac{1}{2}}$ means the square root of a . Similarly $a^{\frac{1}{n}}$ where n is a whole number means $\sqrt[n]{a}$, the n th root of a , and the expression $a^{\frac{m}{n}}$, at first unintelligible, excepting when n divides exactly into m , simply means $\sqrt[n]{a^m}$, the n th root of the m th power of a .

The idea of fractional indices is capable of very important developments. 10^0 is equal to 1, and 10^1 is equal to 10. Further, $\sqrt{10}$ or $10^{\frac{1}{2}}$ is about equal to 3.2. It appears, therefore, that if we select any number between 1 and 10 this is equal, approximately, to some power of 10. 2, for instance, should be some power of 10, and it is easy to see, as a matter of fact that 2 is nearly $10^{\frac{1}{10}}$, for raising 2 and $10^{\frac{3}{10}}$ to the tenth

power we obtain $32 \times 32 = 1024$, and 1000, showing that a fractional power of 10, slightly larger than 0.3, will be equal to 2. Similarly the $\frac{12}{25}$ power of 10 is nearly equal to 3. This leads us to the idea that to all numbers, whole and fractional between 1 and 10, there corresponds some fractional power of 10. These powers, as is generally known, are called the common logarithms of the corresponding numbers. We see easily that if we multiply any number a , between 1 and 10, by 10 this number becomes $a \times 10^1 = 10^m \times 10^1 = 10^{m+1}$ where m is the logarithm of a ; if we multiply by any other power of 10, the logarithm is equal to that of the original number plus this power of 10. We shall return more than once to this important matter of logarithms.

Positive Real Numbers. We have seen that if in a straight line, starting at a point 0 and proceeding indefinitely to the right, equidistant points are located, the distances between these points and the start 0 correspond to the natural numbers. Any number of other points between, say, those corresponding to 1 and 2, can be found by equal subdivision of the length 1 to 2, but these additional points can never make up the whole line. An unlimited further number of points, distant from 0, by amounts representing the irrational numbers can be interpolated between the fractional number points, but even these cannot complete the straight line, because there is a further class of numbers, called transcendental, of a more complicated character than the irrationals, which cannot be expressed arithmetically by addition, subtraction, multiplication, or division of roots of natural numbers. The number π , the sines and tangents of nearly all angles, are examples of this further class, transcendental numbers. The points on the number line representing the fractional, the irrational, and the transcendental numbers together make up the whole line, and the numbers corresponding to all the points on the line make up a totality which is known as that of the positive real numbers. We shall very shortly see the reason for the use of the qualifying terms "positive" and "real."

Negative Numbers. The extension of the conception of number from that of the natural or cardinal numbers used in counting, to include fractional, irrational, and transcendental

numbers has largely removed the restrictions with which the ordinary operations of arithmetic are inherently surrounded. Fractional numbers make the scope of division unlimited, and irrational numbers do the same for the process of evolution. There is one operation of basic arithmetic still restricted, that of subtraction. The sum $3 - 4$ is, so far, meaningless because, as $3 - 3$ is zero or nothing, $3 - 4$ appears to indicate an absurd answer of less than nothing.

As a matter of fact it is easy to find concrete examples of an intelligible meaning of less than nothing. A person of good standing who is allowed to overdraw his current account is less than nothing in credit. Again, if a man walks 2 miles eastward from his house and reduces his eastward movement by walking 1 mile westward, his eastward displacement will be 1 mile, but if he reduces his eastward displacement by walking 3 miles westward he will pass his house, and his eastward displacement will be 1 mile less than nothing. So, in our geometrical conception of real numbers as being distances on a line measured from a starting point 0, the subtraction of one distance from another which is less could be conceived to have an answer if the number line, instead of stopping at the point 0, extended indefinitely in both directions.

In order to develop this idea rationally we conceive ordinary numbers, represented by distances, say to the right from a point 0 on a line, to have a directive character. Thus an ordinary number, 3, means a step of 3 from 0 to the right. Displacements to the right are additive, plus, or positive. A displacement to the left is reckoned subtractive, minus, or negative. According to this directive idea of number the sum $3 - 4$ means a positive number $+3$, on which is superposed, or to which is added a negative number -4 ; the answer to the sum is -1 , a negative number objectively represented by a distance 1 measured to the left of the zero point 0 on the number line. According to this view common subtraction becomes, in our geometrical interpretation of number, equivalent to the addition of a negative number, and the answer will have the character or sign of the number which has the greater magnitude or which represents the greater geometrical length on the number line.

We can easily find a general meaning for subtraction, as the subtraction of a positive number is equivalent to the addition of a number equal in magnitude but opposite in sign or direction, so subtracting a negative number, or a backward, leftward movement, means the addition of a movement to the right. The subtraction of a negative number is therefore equivalent to the adding of a positive number representing the same geometrical length on the number line. This idea is confirmed by the banking illustration: if a man reduces his overdraft his credit at the bank improves.

We can now express our preliminary ideas in symbolical form. The elementary arithmetical operation $a - b$ can be written $+a - b$, meaning, on a positive step a superimpose a negative step b . This is seen at once to be the same as $-b + a$, which means that the positive step is to be superimposed on the negative step. Thus, the operation $a - b$ is commutative, when we consider the $-$ sign to be a description of the nature of the number whose magnitude, direction ignored, is b . Secondly, we have seen that $+a - (-b)$, the subtraction of $-b$ from $+a$ is simply $+a + b$, represented objectively by the superposition of two positive steps along the number line. If the first number or symbol of a mathematical statement like $a - b$ bears no sign, this number is always regarded as positive.

So far our development of an idea of numbers less than nothing, or negative, has been fairly straightforward. If, however, these new numbers are to be of real use to us they must be able, not only to enter into the operations of addition and subtraction, but also into those of multiplication and division, and it is in these latter fields that real difficulties arise.

Let us consider a very obvious idea, that of the multiplication of a directed step on the number line, in either direction, in the sense of repeated addition. 3 times a positive step of 2, or 3 times $+2$ evidently gives $+6$ as an answer, and 3 times -2 gives -6 . We need not limit multiplication to whole numbers for $3\frac{1}{2}$ times a negative step -2 would evidently give -7 . The difficulty here is, however, that we are really considering three kinds of numbers, positive and negative, and a multiplying or magnifying number which inherently possesses nothing of the nature of direction. -2×3 , a step -2 ,

increased threefold is intelligible, but the operation -2 multiplied by $+3$, or $(-2) \times (+3)$ is inherently meaningless, because a multiplier is devoid of direction. Let us agree, however, that the $+$ sign attached to a multiplying number means multiplication in the usual sense, and that $(-2) \times (+3)$ means 3 times -2 , or -6 ; what meaning are we to attach to the commuted operation $(+3) \times (-2)$, that is -2 times 3? Here the operation is inherently meaningless, but, if the commutative law of multiplication is to hold, this should give the same result as $(-2) \times 3$ or -6 . If we agree to this meaning, we arrive at the conclusion that the multiplication of two numbers of different character or sign always gives a negative answer. There is, however, a further kind of multiplication to be considered, even more unintelligible inherently, that of the multiplication of two negative numbers, and, if this is to give an answer, this answer can only be positive, for as $(+a) \times (-b)$ gives $-ab$, $(-a) \times (-b)$ must give an answer of different sign. Let us see if this rule will work in an arithmetical example. Consider the sum $8 - 2(3 - 2)$, which means subtract from 8 twice the difference of 3 and 2; the answer is 6. If, however, we consider the multiplier 2 to be a negative number, the sum tells us to superpose on the positive step 8 a negative step $(-2) \times (+3 - 2)$ or two steps, the one $(-2) \times (+3)$ and the other $(-2) \times (-2)$. The first of these component steps is according to our rules a negative step -6 , the next $(-2) \times (-2)$ is $+4$. Thus, according to the rules or conventions for multiplication by negative numbers, our arithmetical sum means geometrically a $+$ step 8, a $-$ step 6, and a $+$ step 4, giving an answer $+6$. This may all seem obvious, but it is of the greatest importance clearly to realise what assumptions we have made, or rather, what conventional meanings we have assigned to operations which inherently are meaningless. In the first place we have assigned a directive character to number, and have thereby arrived at a set of entirely new numbers having a directive character negative, which is opposite to that, positive, of the numbers previously considered. Further, and what is most important to realise, we have given a kind of directive character to multiplication, and assumed that this operation by a negative

number, not only changes the magnitude, but also the directive character or sign of a number multiplied.

The extension of our ideas of division leads to no difficulties once these ideas of multiplication are grasped, for as $\frac{-a}{(b)}$ is plainly $-\left(\frac{a}{b}\right)$, so $\frac{a}{-b}$ must give an answer which, multiplied by $-b$, is simply a . This answer must be $-\left(\frac{a}{b}\right)$ or $\frac{-a}{b}$, for $\frac{(-a) \times (-b)}{b}$ is $\frac{+ab}{b}$ or simply a . $\frac{-a}{-b}$ also is plainly $\frac{a}{b}$. Thus the rules for division, as regards the character or sign of the answer, are exactly the same as those for multiplication. We should note carefully that these rules cannot be said to have been proved; we cannot really prove something which is inherently meaningless; rather they have been tested and found to work, or to lead to consistent results.

Keeping to the geometrical interpretation of number we see that to every ordinary or positive number there corresponds a similar negative number, representing a reverse step from the 0 of the number line of exactly similar length. Every point on a line unlimited in extent in both directions with a reference point 0 therefore represents some number, positive or negative; the numbers making up this unlimited totality are called the real numbers.

Negative Indices. It is relatively easy to find a meaning for the inherently absurd operation of raising a number to a negative power, for as dividing a^2 by a reduces the index by 1 and gives $a^1 = a$, and further division by a reduces the index to 0 and gives $a^0 = 1$, so a third division by a reduces the index by 1 from 0 to -1 , and $\frac{1}{a}$ is found to be equal to a^{-1} . Further, dividing a^0 or 1 by a^2 reduces the index from 0 to -2 so that $\frac{1}{a^2}$ is the same as a^{-2} . Generally, $a^{-m} = \frac{1}{a^m}$, and it can be shown that, generally, $a^{-\frac{m}{n}}$ is equal to $\frac{1}{a^{\frac{m}{n}}}$, when a is a positive real number.

The negative index notation is very useful for expressing small decimal fractions. As 0.000006 is 6 millionths, and as 1 millionth is 1 divided by 1 million, or 10^{-6} , so 0.000006 can be expressed as 6×10^{-6} . The negative index is equal to the number of figures in the decimal fraction, and the actual figures of this fraction become a whole number. For another illustration 0.00361 is equal to 361×10^{-5} . The index notation is convenient for the multiplication and division of very large or very small numbers. Thus 0.000012×0.004 is $12 \times 10^{-6} \times 4 \times 10^{-3} = 48 \times 10^{-9}$ by addition of indices of 10. Again $6,240,000 \div 0.00000072$ is $624 \times 10^4 \div 72 \times 10^{-8}$, or $10^{12} \times 624 \div 72$, by subtraction of the index -8 from the index 4.

We have seen that numbers between 1 and 10 can be expressed approximately as fractional powers of 10. Thus $2 = 10^{0.3}$ nearly. Further, $2 \times 10^n = 10^{n+0.3}$, and multiplying 2 by a power of ten adds n to this power or logarithm, so that the logarithm of 2000 is about 3.3. Since 0.2 is 2×10^{-1} , $0.2 = 10^{-1} \times 10^{0.3} = 10^{0.3-1}$. The logarithms of numbers less than 1 are negative because the logarithm of 1 is 0. It is customary to express negative logarithms as the sum of a negative whole number, which represents the order of the number, and a fractional part, always positive which depends only upon the actual figures composing the number. This latter part is called the mantissa; the variable part, depending upon the position of the decimal point in the set of figures forming the number is called the characteristic. The logarithm of 2000 is approximately 3.3, that of 0.02 or 2×10^{-2} is expressed as $\bar{2}.3$, this symbol meaning $-2 + 0.3$. A positive characteristic of a number over 10 is one less than the number of figures in it, the negative characteristic of a positive number less than 1 is the same as the number of figures or digits of which it is composed.

It is evident that the larger the index of a power of 10 the larger is the power, and the larger is the negative index the larger is the denominator of the fraction representing the power. No negative number can have a logarithm in the sense of the word so far considered, that is, no power of 10, the index of which is a real number, positive or negative, can be negative.

Involution and Evolution of Negative Numbers. The involution of negative numbers in the sense of repeated multiplication gives rise to no difficulties. The square of $-a$, or $(-a)^2$ is $(-a) \times (-a)$ which, as we have seen, must have a positive sign $(-a)^2$ is therefore a^2 . Similarly $(-a)^2 \times (-a)$ or $(-a)^3$ is $a^2 \times (-a)$ or $-a^3$. It appears, therefore, that even powers of a negative number are positive and odd powers negative.

No general meaning can at present be assigned to the evolution of negative numbers, or the raising of them to fractional powers. For it is impossible to discover a real number which, when squared, gives a negative answer. The square root of -4 , for instance, is certainly not $+2$, and it cannot be -2 , for -2×-2 is $+4$. We can, however, assign a cube root to a number like -8 , for if -2 is cubed we obtain -8 as the answer. We conclude then that as odd powers of negative numbers are negative, so negative numbers have negative roots of odd orders, but no numbers of the kind we have considered can stand for any even root of a negative number. Our final task in this chapter on number will be to discover a new class of numbers, among which are some having the property that their squares are negative real numbers.

Extended Idea of the Evolution of Positive Numbers. The reader will probably have remarked that as the square of a negative number, such as -2 is positive, the number 4 appears to have two square roots $+2$ and -2 . This is often expressed symbolically as $\sqrt{4} = \pm 2$. The cube number 8 , however, has only cube root 2 ; while 16 has two fourth roots, $+2$ and -2 . We shall shortly see that, with an extended conception of number, not only are we able to assign a meaning to the square root of a negative number, but that any number has a number of roots equal to the order of the root, for example, that there are two cube roots of 8 each different from 2 , and each different from the other. Similarly 16 will be found to have two fourth roots additional to $+2$ and -2 .

Symbolic Representation of Directed Steps. Our development of an unlimited domain of real numbers, positive and negative, has been based largely upon the objective geo-

metrical conception of directed steps in an unlimited straight line with a reference or zero point. The reader may have asked himself the question, can any number significance be allotted to steps from the reference point outside the real number line? We have already dealt in a rudimentary way with such steps in Chapter III. The treatment was there geometrical; to interpret these steps as having a number significance, such as has been assigned to steps in the number line, we must be able to show that they are amenable to the ordinary arithmetical operations of addition, subtraction, multiplication, and division, with perhaps the addition of conventions which are supplementary to, but which do not contradict the fundamental rules.

In Fig. 21 we have represented a step Z_1 outside the number line X_1OX , of which the points in the section OX represent

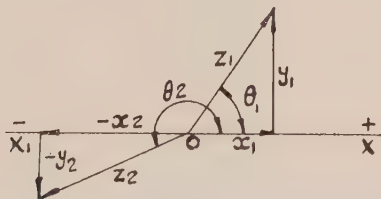


FIG. 21.

positive real numbers, and those in OX_1 negative numbers. This step Z can, as we have seen in Chapter III, be made by the superposition of two component steps, the one x_1 in the number line and the other y_1 perpendicular to it. The recognition of negative numbers enables us to adopt a much more concise notation for directed steps than was possible in Chapter III. For, as we reckon steps in the OX direction positive and those in the OX direction negative, so we can consider upward perpendicular steps positive and downward perpendicular steps negative. The step like Z_2 is therefore compounded of a horizontal step $-x_2$ and a perpendicular step $-y_2$. If we consider horizontal steps normal and vertical steps abnormal we might indicate the step Z_1 as $x_1 + y_1$ (perpendicular). It would evidently be advantageous to decide

on a characteristic and concise mark which, attached to a length symbol, would indicate that this referred to a perpendicular step. Two marks of this kind are used in mathematics, these are the letters i and j , placed before the length symbol. Thus the vertical compound step of Z_1 can be indicated iy , or jy , the i and j at present being merely marks, something like the δ used on p. 60. It is unfortunate that two symbols of this kind are used, i in pure mathematics, and j in electrical engineering; as this book is written primarily for engineers we shall use the j symbol. $x+jy$ therefore means a step x length units in the positive direction along the

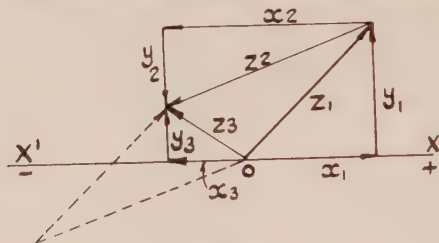


FIG. 22.

number line, on which is superposed a step y upwards and perpendicular to the number line.

We can soon satisfy ourselves by a brief study of Fig. 22 that the combination of two steps Z_1 and Z_2 can be represented by this kind of notation. The resultant of these two steps is a third step Z_3 , of which the horizontal component is a step equal to $x_1 - x_2$, or the sum given by adding to the positive steps $+x_1$ the negative step $-x_2$, and the perpendicular component is given by adding to the perpendicular step upwards $+y$, a perpendicular step downwards $-y_2$. As therefore we can symbolically represent the step Z_1 in the form $Z_1 = x_1 + jy_1$, and the step Z_2 in the form $Z_2 = -x_2 - jy_2$, so we can put $Z_3 = Z_1 + Z_2 = x_1 - x_2 + jy_1 - jy_2$, or as the net perpendicular step in Z_3 is $y_1 - y_2$, this last statement can be abbreviated $Z_3 = x_1 - x_2 + j(y_1 - y_2)$. It is easy to see by the dotted line construction in Fig. 22 that this addition of vectors in a symbolic way is commutative, and that we should have obtained the same result by adding Z_1 to Z_2 .

The subtraction of vector steps can be interpreted by the ordinary rules for the subtraction of real number steps. To subtract a real number b from a we change the sign of b and add this changed number to a . If b is positive, adding its reversal diminishes a , if b is negative, adding its reversal $+b$, increases a . So, to subtract from a step Z_1 a second step Z_2 we reverse the direction or sense of Z_2 and add this reversed step to Z_1 . The reversed step Z_2 of Fig. 22 becomes $-Z_2 = x_2 + jy_2$ so that $Z_1 - Z_2$ is expressed symbolically as $x_1 + x_2 + j(y_1 + y_2)$, and this symbolic expression is shown geometrically in Fig. 23. We see that subtraction in this sense is commutative, for $Z_1 - Z_2$ is plainly the same as $-Z_2 + Z_1$.

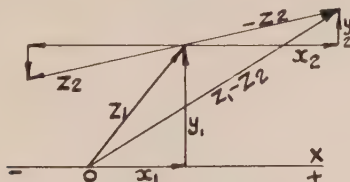


FIG. 23.

We had better now review briefly what we have done in the matter of the identification of geometrically represented steps with a symbolic notation. The symbol Z applied to a directed step and standing for the superposition of two component steps contains implicitly two items of information, first the measure of the length of the step r which we have seen on p. 53 is equal to the square root of $x^2 + y^2$, or to $\sqrt{(x^2 + y^2)}$, and secondly the angular inclination of the steps to the number line which evidently depends upon the ratio of y to x . The measure of the length of a step r is always reckoned a positive number, whatever its direction may be. The direction or inclination of a step is reckoned by the angular amount θ the line representing the step has been turned from the OX line of positive real numbers. This angular turning is reckoned positive if anti-clockwise. In the case of Z_2 in Fig. 21, the angular turning or inclination θ is a positive amount over half a complete turn or over two right-angles.

By considering the step Z_1 of Fig. 21, and by recollecting the definitions of the trigonometrical ratios of an angle given on p. 44, we see that as $\cos \theta_1$ is equal to $\frac{x}{r}$, and $\sin \theta_1$ to $\frac{y}{r}$, the x component step is $r \cos \theta$, and the y component step $r \sin \theta$, so that the symbolic expression $x + jy$ can be written $r \cos \theta + j(r \sin \theta)$. So far the trigonometrical ratios have been considered only in reference to the angles of a right-angled triangle, that is for values of θ not greater than 1 right-angle. These ratios are defined generally in such a way that the symbolic statement of a directed step can always be represented in the form $r \cos \theta + j(r \sin \theta)$, θ being the anti-clockwise angular turning of the step relative to the OX positive number line. Thus considering Z_2 in Fig. 21, as x and y are both negative, so are $\cos \theta$ and $\sin \theta$, for r is always positive.

Symbolic expressions for directed steps can be added and subtracted like ordinary numbers, in fact these expressions show that the steps are composed each of two distinct kinds of components, one of which is distinguished by the mark " j ," and that the combination or addition of the steps means adding their similar components. If however these symbolic expressions are to have the properties of numbers it must be possible in some way to multiply and divide them, and these operations are generally without apparent meaning.

We have already, however, considered one case of the multiplication of directed steps in the number line, and we found that there is no inconsistency in regarding the multiplication of $+2$ by -3 as the same as the obvious process of doubling the step -3 . The doubling of a step $x + jy$, of length r , and inclination θ is fairly obvious, the result is a step in the same direction of double the length, which can be represented $2x + j(2y)$, for each of the component steps will have been doubled. To multiply a step $+2$ in the number line by the inclined step $x + jy$ is really meaningless; let us however agree that this means the same as multiplying $x + jy$ by 2. The result of the multiplication of $+2$ by a directed step of length r and inclination θ is to alter the length of $+2$ to $2r$, and to turn it through an angle θ ; in short, to increase its inclination from zero to θ . If we generalise this idea and say that the

multiplication of two directed steps, the one characterised as r_1, θ_1 , and the other r_2, θ_2 , gives a step of length $r_1 r_2$ and angular inclination $\theta_1 + \theta_2$, we shall have the rule for the multiplication of vectors or directed steps. We can put this rule in the concise form : to multiply vectors, multiply lengths and add angles.

We must now see how this rule fits in with the symbolical or algebraic statement of a vector, for if it is to work successfully we must be able to obtain the product of $x_1 + jy_1$ by $x_2 + jy_2$ by multiplying them according to the rules of generalised arithmetic or algebra. Let us consider the step Z of Fig. 24 compounded of direct horizontal and perpendicular

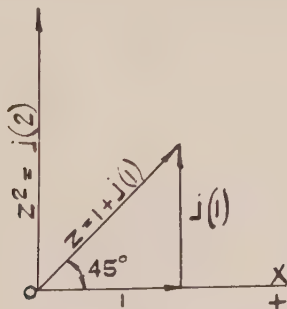


FIG. 24.

steps, 1 and $j(1)$, and let us multiply this step by itself. We see at once that the length of Z is $\sqrt{2}$ and that its inclination is an angle of 45 degrees. According to the rule for vector multiplication the product of Z by Z should give a vector of length $\sqrt{2} \times \sqrt{2}$ or 2 and inclined $45 + 45$ or 90 degrees to the OX positive number line. This is a purely perpendicular step indicated $j(2)$. Now let us obtain the product by algebra, multiplying $\{1 + j(1)\}$ by $\{1 + j(1)\}$. We obtain $1 + j(1) + [j(1)] \times [j(1)]$ or, combining the two similar terms, $1 + j(2) + [j(1)] \times [j(1)]$. This does not look right at first sight, but before we decide this point let us examine the term $[j(1)] \times [j(1)]$. If instead of considering j as a mark we regard it as some kind of number, this term is $j \times j$, and now we see that our answer $1 + j \cdot 2 + j \times j$ will agree with the answer obtained by vector multiplication if $j \times j$ is equal to -1 , for this -1

will cancel the 1 and we shall be left with $j(2)$. Let us go further and multiply $j(2)$ by itself. Vectorially this should give a step 2×2 or 4 in length inclined at two right-angles to the OX direction. This is a step of 4 in the negative number direction or -4 . Treating $[j(2)] \times [j(2)]$ algebraically this is $j \times j \times 4$, and if $j \times j$ is -1 , the answer is -4 . Thus we see that if j is regarded, not as a mark, but as some kind of number such that its square is -1 , algebraic multiplication of vectors agrees with the rule, multiply lengths and add angles at any rate in the simple cases considered.

Having formulated the multiplication rule, the rule for division follows easily, for if Z_3 has a length r_3 and an inclination θ_3 , then dividing it by Z_2 , of length r_2 , inclination θ_2 gives a length $\frac{r_3}{r_2}$, and an inclination $\theta_3 - \theta_2$, for multiplying this vector by the divisor r_2 we evidently arrive at the original vector Z_3 . The rule for division is therefore: divide lengths and subtract angles.

Imaginary and Complex Numbers. We now see that we have filled up the original gap in our complete number scheme, that of the absence of numbers with negative squares, by the invention of a new kind of number j , having the basic property that its square is -1 . j is evidently, in the geometrical sense, a unit step in a direction perpendicular to the real number line, for squaring it or doubling its inclination of one right-angle makes it into a unit step in the negative number line or turns it into -1 . Any negative number, $-a$, being $-1 \times (+a)$, has therefore the square root $j \times \sqrt{a}$.

A number represented symbolically ja , and geometrically by a perpendicular step upwards of length a , is called an imaginary number, $-ja$ is a reverse step that exactly cancels ja and is therefore a perpendicular step downwards. Multiplying a real number by j , or, in our original sense, attaching the j mark to it, means multiplying its length by 1, and adding to its inclination an angle of 90 degrees.

There has been a good deal of mystery surrounding this new number j , often called the square root of -1 , and perhaps on this account the name imaginary has been given to it. j is really no more imaginary than an irrational surd number,

or a negative number. As no combination of natural numbers can possibly represent the new kind of number indicated by $\sqrt{2}$, and as no answer is obtainable to the sum $3 - 4$ without the invention of a new class of numbers of which -1 is the type, so j is a type of a further new class of numbers having the inherent property that their squares are negative. The geometrical interpretation of j is as objective and real as the geometrical interpretation of -1 .

A number like $x + jy$ composed of a real part x and an imaginary part jy , and represented geometrically by an inclined vector, is called a complex number. Complex numbers obey the laws of elementary algebra as do real numbers, if in calculating with them $j \times j$ is always put equal to -1 . As the operation of multiplying complex numbers comprise multiplication of lengths and additional of angles, in the geometrical sense, this is evidently commutative, in that $Z_1 \times Z_2$ gives the same result as $Z_2 \times Z_1$.

Numbers like $x + jy$ and $x - jy$ are of equal length geometrically and of opposite inclinations. Multiplied together the inclinations ought to cancel and give a real product. Algebraically the product is $x^2 + (jy)^2$ or, as $j^2 = -1$, $x^2 + y^2$, a real answer, equal in magnitude to the square of the length of each of the factors. Two such numbers are said to be conjugate, and it is evident that a fraction like $\frac{a}{x + jy}$ can be made to have a real denominator by multiplication of both numerator and denominator by $x - jy$, giving a result $\frac{a(x - jy)}{x^2 + y^2}$.

We have seen that $x + jy$ can be written in the form $r \cos \theta + j(r \sin \theta)$, and as j is now found to be a number this can be abbreviated to $r(\cos \theta + j \sin \theta)$. This expression for a complex number shows its two geometrical characteristics explicitly; $r = \sqrt{(x^2 + y^2)}$ is geometrically the length of the directed step, and is called the modulus. θ , an angle whose tangent is $\frac{y}{x}$, is called the argument. The portion $(\cos \theta + j \sin \theta)$ of the expression is sometimes abbreviated to the symbol $\cos \theta$. If $r(\cos \theta + j \sin \theta)$ is a correct algebraic statement of a complex number, it ought to be amenable to

rules of arithmetic. We see that the addition of two such expressions does not lead to any apparent useful results, but we may profitably see the result of multiplication. Consider two numbers each with unit modulus ($\cos \theta_2 + j \sin \theta_1$) and ($\cos \theta_2 + j \sin \theta_2$). In the geometrical sense these two unit steps should when multiplied give a third unit step with an inclination of $(\theta_1 + \theta_2)$, that is, one represented by $\cos (\theta_1 + \theta_2) + j \sin (\theta_1 + \theta_2)$. The two component steps multiplied algebraically give

$$(\cos \theta_1 + j \sin \theta_1)(\cos \theta_2 + j \sin \theta_2) \\ = \cos \theta_1 \cos \theta_2 - \sin \theta_1 \sin \theta_2 + j(\sin \theta_1 \cos \theta_2 + \cos \theta_1 \sin \theta_2)$$

since $j \sin \theta_1 \times j \sin \theta_2$ gives $-\sin \theta_1 \sin \theta_2$, j^2 being replaced by -1 .

The two forms of the product each contain a real and an imaginary part, and if they are equivalent, the real part of the one must be the equivalent of the real part of the other, and, similarly with the imaginary parts. If this is the case,

$$\begin{aligned} \cos (\theta_1 + \theta_2) &= \cos \theta_1 \cos \theta_2 - \sin \theta_1 \sin \theta_2 \\ \text{and} \quad \sin (\theta_1 + \theta_2) &= \sin \theta_1 \cos \theta_2 + \cos \theta_1 \sin \theta_2 \end{aligned}$$

which statements are proved geometrically in the textbooks on trigonometry, and thus confirm the consistency of the rule for the multiplication of complex numbers which represent directed steps.

Involution and Evolution of Complex Numbers. The multiplication rule for complex numbers leads at once to rules for involution and evolution, for raising a number having a modulus r and an argument θ to the n th power, by repeated multiplication, means raise r to the n th power and multiply θ by n . The n th power is thus a complex number with a modulus r^n and an argument $n\theta$. If we use this rule to square $(\cos \theta + j \sin \theta)$ we obtain $(\cos 2\theta + j \sin 2\theta)$. Squaring by algebra we obtain

$$(\cos \theta + j \sin \theta)^2 = \cos^2 \theta + j^2 \times \sin^2 \theta + 2j \cos \theta \sin \theta$$

and as real and imaginary parts of the two expressions must be equal

$$\cos^2 \theta = \cos^2 \theta + j^2 \times \sin^2 \theta = \cos^2 \theta - \sin^2 \theta.$$

$$\text{and} \quad \sin 2\theta = 2 \cos \theta \sin \theta.$$

It is evident that proceeding in this way we could obtain equivalent expressions for $\cos n\theta$ and $\sin n\theta$, by raising $(\cos \theta + j \sin \theta)$ to the n th power by algebraic calculation.

Evolution of complex numbers can similarly be readily interpreted, for as the n th root of a number of modulus r and argument θ , raised to the n th power must merely reproduce the number, the root is simply a number of modulus $\sqrt[n]{r}$ and of argument $\frac{\theta}{n}$.

This idea of involution, moreover, explains why a number like 4 has two square roots, for as a positive number 4 can be rotated through one complete turn or four right-angles to produce itself, this identical and rotated number has a square root with half the rotation, that is, 2 in the negative number line or -2 . Similarly as $+1$ is the same after one or two complete revolutions, so the cube root of $+1$ once completely turned is a complex number represented by ω in Fig. 25, of unit

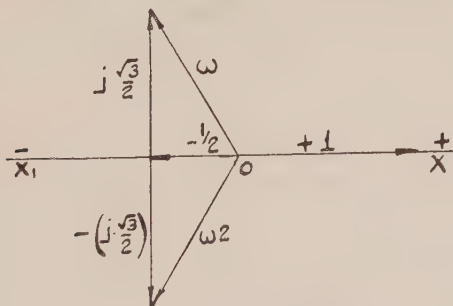


FIG. 25.

modulus, with an argument $\frac{1}{3}$ of a turn or 120 degrees. This number it can be seen from the geometry of the diagram to be the same as $-\frac{1}{2} + j\frac{\sqrt{3}}{2}$. Similarly $+1$ rotated through two complete turns has a cube root represented by ω^2 , with an argument of $\frac{2}{3}$ of a turn, and this is seen to be the same as $-\frac{1}{2} - j\frac{\sqrt{3}}{2}$. ω^2 is evidently the square of ω for it has double the argument or angle relative to the positive number 1,

$\omega^2 \times \omega$ ought then to equal 1. We can easily see that this is so for the algebraic multiplication gives $(-\frac{1}{2})^2 - \left(j \cdot \frac{\sqrt{3}}{2}\right)^2$ or $\frac{1}{4} - j^2 \times \frac{3}{4} = \frac{1}{4} + \frac{3}{4} = 1$.

It can similarly be shown that +1 or unity has n n th roots. The reader should satisfy himself that the four 4th roots of 1 are 1, -1, j and $-j$. The n roots are evidently given geometrically by unit radii of a circle which divide the circumference into n equal parts, one of these radii corresponding in direction with the positive real number line. These roots can be expressed in the form $x + jy$ only when the corresponding equal subdivision of the circumference of a circle is possible by common geometry (see p. 47).

Completeness of the Number System. All the operations of arithmetic, addition, subtraction, multiplication, division, involution, and evolution can be carried out without restriction on numbers comprehended within the complete system that we have developed, real, imaginary, and their combination, complex numbers, and it can be shown that all combinations of these operations are possible, and can lead to no requirement of a further new kind of number for their rational interpretation.

The reader may rightly inquire whether, as we have devised new numbers represented geometrically by lines lying outside the line of real numbers, we can go further and devise further new numbers lying outside the plane of the complex numbers; in other words, whether we can have a further kind of unit, additional to 1 and j , and geometrically at right-angles to both. Systems of numbers have been devised on the basis of this idea, and such numbers are called quaternions. The algebra of quaternions is, however, much more complicated than that of real and complex numbers because the commutative law of multiplication does not hold with them; in other words, a product depends upon which of the two factors is the multiplier. Quaternions does not enter into the elementary applications of mathematics to engineering, and we shall therefore do no more than make this brief mention of them.

CHAPTER V

FUNCTIONS

Elementary Considerations. We have hitherto considered the mathematical idea of magnitude, as expressed by a number, in a stationary sense. We have now to consider the idea of variable magnitude. In the material world change is normal, and the development of physical science has been largely the investigation of the quantitative laws which, by experiment and experience, have been found to express the variation of a secondary magnitude or effect produced by a variation of a primary magnitude or cause. A mathematical expression that describes the connection of one magnitude with another is called a function. The primary magnitude, which may objectively be thought of as measuring a cause, is called the independent variable, and it is, in pure mathematics, denoted by the letter symbol x . The secondary quantity which may be thought of as measuring an effect is called the dependent variable, and it is denoted by a letter symbol from the end of the alphabet. If a quantity y is a function of x , then, assigning any value to x , the corresponding value of y can be calculated. A relation of this kind, between x and y , is denoted symbolically as $y=f(x)$, which is a shorthand way of writing, what is read " y is a function of x ," that the value or values of y are completely determined by the value of x . A dependent variable may be a function of two independent variables; this would be expressed symbolically as $u=f(x, y)$, and this means that u is affected by changes of either or both x and y , but values of x and y being assigned, the value or values of u can be found. Functional relationship between y and x is also expressed symbolically as $y=\phi(x)$, or $y=\psi(x)$ when the letter symbol " f " is required for other purposes. All the foregoing is a matter of definition.

The simplest example of functional relationship is what is often called direct proportion. Thus, in the use of electrical energy under a flat-rate tariff, the amount payable for the supply is directly proportional to the consumption, and if y is

the amount payable, a the rate per unit, and x the consumption, $y = ax$. Here of course mathematical symbolism is mere verbosity, from a practical point of view, but it is used to illustrate the simplest of all functions that can be expressed in algebraic symbolism. With a stipulated value of the rate a , then the amount of the bill can be determined for any value of the unit consumption.

A slightly more complicated example of a function is given by the bill for electricity under a two-part tariff, when the amount payable is a fixed quarterly charge plus a charge proportional to the consumption. This could be expressed algebraically as $y = b + ax$, where b is the fixed charge. We observe that if a is pence per unit, b must be the amount of the fixed charge in pence, and y will be the total bill in pence. Other simple examples of functions are the formulae found in engineering hand-books of reference.

The quantity a in the function $y = ax$ is deemed to remain constant as x varies. Quantities of this kind which occur in mathematical functions are called constants and are indicated by letters at the beginning of the alphabet. Sometimes it is necessary to study a group of functions, all expressed by the same symbolism or formula, but in which one or more of the constants take different values. In such a case the constants are sometimes called parameters.

Functional relation between two variable quantities can exist and be perfectly explicit, although this relation cannot be expressed in algebraic symbolism. We can, for instance, invent a function of x which takes the value 2 when x is a whole number, but which is 1 when x has any other value. y in this case is a function of x , because assigning any value to x , say 2.6 y is known to be 1, and when $x = 5$, $y = 2$; but this functional relationship cannot be expressed explicitly in any but the verbal form.

When the value or values of a dependent variable y are expressed mathematically in the form $y = f(x)$ so that, given x , the determination of y is merely a matter of calculation, the function is called explicit. $y = b + ax$ is an explicit function of x . An implicit function is a mathematical expression containing both x and y , but in such a form that it is not

necessarily possible to write a formula connecting y and x . An implicit function is written $f(x, y)=0$. Thus, by changing $y=b+ax$ to $y-b-ax=0$ we obtain an implicit function of y . Any explicit function can be put into the implicit form, but the converse is usually impossible. Thus $y^5-y-x=0$ is an implicit function of x , because, assigning a value to x fixes the values of y which make the mathematical statement true, but this cannot be converted to the form $y=f(x)$ in which the value of y in terms of x is expressed explicitly.

Engineering Formulae. In the practical use of the simple functional relationships expressed by what are often called formulae, it is sometimes necessary to invert the function in the sense of finding the value of the independent variable that corresponds to an assigned value of the dependent variable. Thus, consider the formula $s=\frac{1}{2}gt^2$ which gives the space s passed over by a body falling from rest during a time t , g being a gravitational constant. Here s is a function of t^2 , and knowing g , the value of s is readily calculated by simple arithmetic for any given value of t . We can readily alter this formula to give t as an explicit function of s . The statement $s=\frac{1}{2}gt^2$ is a simple example of what is called an equation, and it is almost self-evident that the two equal quantities s and $\frac{1}{2}gt^2$ can each be divided or multiplied by the same quantity and still retain their equality. So, using 2 as a multiplier, we obtain $2s=gt^2$, and using g as a divisor, $\frac{2s}{g}=t^2$, whence

$t=\sqrt{\left(\frac{2s}{g}\right)}$, an explicit functional relationship between s and t .

Again, $t=2\pi\sqrt{\left(\frac{l}{g}\right)}$ gives the time of oscillation of a simple pendulum as an explicit function of the length l , 2π being a number, and g the gravitational constant. As the squares of equal quantities are equal, $t^2=\frac{4\pi^2l}{g}$, so that we can, in the manner shown in the preceding paragraph, obtain the equivalent equation $l=\frac{t^2g}{4\pi^2}$, in which, assigning a value to the time of oscillation t , the corresponding length l is explicitly given.

As a further example of this kind of thing, consider the function $y = b + ax$, which we have seen expresses the charge for an electricity supply under a two-part tariff, with a fixed quarterly charge. It is not often that we should require to know the consumption corresponding to a stipulated value of the total charge, but this could readily be expressed algebraically. For if from each of the equal quantities y and $b + ax$ we deduct the same quantity b , we obtain $y - b = ax$, from which $x = \frac{y - b}{a}$.

This kind of algebraic manipulation is mainly a matter of the technique of calculation, although, anybody desiring to understand mathematics, should be so familiar with it that mental assent can be given intuitively. The rules for this kind of manipulation are sometimes given in the following form: any quantity can be transferred from one side of an equation to another by changing its sign; a multiplying factor of the whole of one side of an equation can be removed and made a divisor of the whole of the other side, and vice versa. The qualification in the last rule must be noted. Thus, if we had a formula $y = a + \frac{x}{b}$, we could not transfer the divisor b , as the equation stands, we should have to make it a divisor of the whole of the right-hand side by changing the equation to $y = \frac{ab}{b} + \frac{x}{b}$ or $y = \frac{ab + x}{b}$. Then we could transfer the divisor b and obtain $by = ab + x$, and, finally, $x = by - ab$.

This manipulation of simple formulae is a very simple example of the much wider matter of the solution of equations, and it is only in a very limited number of cases that the transformation can be made. We shall return to this subject later.

Homogenous Functions—Dimensions. We drew attention, on p. 22, to the fact that if an algebraic expression like $(a + b)$ is involved, and the power is expanded, the sums of the powers of a and b in the terms of the expansion are all equal to the index of the power of $(a + b)$. An expression of this kind, in which the sum of the indices of the powers of all the terms is the same, was there called homogenous. A functional expression which represents the relation between

objective or physical magnitudes always has a peculiar kind of homogeneity, in that each separate term of the function must stand for the same kind of quantity.

To illustrate this important point in a very elementary way, consider the expression $(a+b)^2$ as standing for the area of a square the side of which is the sum of two lines of respective lengths a and b . The algebraic expansion of $(a+b)^2$ is $a^2 + 2ab + b^2$, and in this expansion each term must represent the same kind of quantity as does $(a+b)^2$, that is, an area. The diagram of Fig. 26 shows clearly that each of these terms

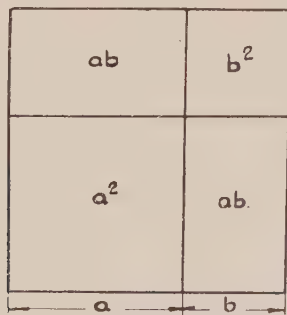


FIG. 26.

stands for an actual area : a^2 and b^2 respectively for the squares having sides a and b , and $2ab$ for the two rectangles of which neighbouring sides of each are a and b .

Again, suppose we encounter a formula like $y = ax^2$, where y and x are both geometrical lengths of lines, then we can be sure that a does not represent a length. But if the formula is given as $y = \frac{x^2}{b}$, then b is a length, for the formula means that the square on x is equal in area to the rectangle of sides y and b .

Any physical magnitude has two distinct characteristics, first the number representing it, which is the ratio of the magnitude to that of a unit of similar character, and, secondly, the nature of the unity. The simplest of all physical magnitudes are length, mass, and time, and the manner in which

these basic quantities are connected with a more complex quantity is expressed by the dimensions of the quantity. Units of measurement of physical quantities other than those of length, mass, and time are called derived units.

Area and volume are measured respectively in units which are the square and the cube of a length unit. This is expressed by saying that the dimensions of area are L^2 and those of volume L^3 . The reason for this kind of notation can be illustrated in a simple way. If lengths are measured in feet, then areas will be measured in square feet, but if lengths are measured in inches, the measure of area in square inches will be 12^2 times the square feet measure, because a foot contains 12 inches, and the dimensions of area are L^2 .

Let us consider another kind of geometrical magnitude, that of angle. Angle is measured as the ratio of an arc to a radius, that is, the ratio of a length to a length. Its dimensions are therefore length per length, $L^1 \div L^1$ or L^0 . Angle is therefore said to be of zero dimensions, otherwise it is said to be expressed by a pure number. This means that, whereas the measure of, say, area depends upon the length unit employed, the measure of angle is independent of this, and is the same, whether arc and radius are both measured in feet or centimetres.

Other important class of derived units are those depending upon length and time. Speed is measured by space divided by time. Its dimensions are therefore said to be, $\frac{L}{T}$ or LT^{-1} , using the negative index notation. Acceleration is measured by change of speed divided by time: its dimensions are therefore those of speed divided by time or $\frac{L}{T^2}$, or LT^{-2} .

Dynamical units depend upon mass as well as length and time. Force is measured by mass multiplied by acceleration, and its dimensions are MLT^{-2} . Work, measured by force multiplied by distance, has therefore the dimensions $MLT^{-2} \times L$ or ML^2T^{-2} .

To realise the application of all this, let us consider the well-known dynamical formula. Kinetic energy $= \frac{1}{2}mv^2$, m standing for mass and v for speed. The dimensions of kinetic energy are the same as those of work, namely, ML^2T^{-2} . The

dimensions of mv^2 are $M \times (LT^{-1})^2$, LT^{-1} being the dimensions of speed, this is equal to ML^2T^{-2} . As the dimensions of kinetic energy and mv^2 are identical, this shows that the $\frac{1}{2}$ in the formula is a pure number, independent of the units in which energy, m and v are measured.

Now consider the formula $s = 16 \cdot 1 t^2$ which might be said to give the space passed over by a body falling from rest, after a time t . The dimensions of the left-hand side of the equation are L , those of the left T^2 multiplied by the unknown dimensions of $16 \cdot 1$. As L and T^2 are different, $16 \cdot 1$ cannot be a pure number, and it must represent some physical quantity the magnitude of which depends upon the units employed. We see at once that the dimensions of the $16 \cdot 1$ are LT^{-2} . These are those of an acceleration, and, as a matter of fact, we know that $16 \cdot 1$ is half the gravitational acceleration in the British Engineers' system of dynamical units.

Another well-known formula that may be tested for dimensional homogeneity is that which gives the periodic time of a simple pendulum, $t = 2\pi \sqrt{\left(\frac{l}{g}\right)}$. The dimensions of $\frac{l}{g}$ are those

of length divided by acceleration or $L \div LT^{-2}$ or T^2 . $\sqrt{\frac{l}{g}}$ therefore has the dimensions T , the same as t . 2π is therefore a pure number, and the formula holds for any system of units, provided of course that these are consistent throughout.

Algebraic Functions. After this short digression into the practical characteristics of simple functions relating to physical magnitudes, we proceed to deal further with the mathematical idea of function and what underlies it.

When the relation of the dependent to the independent variable can be expressed by a mathematical statement of a limited number of the operations of arithmetic, the statement is said to be an algebraic function. Thus $y = \frac{x + ax^3}{c + \sqrt{x}}$ is an algebraic function of x , because assigning a value to x , the corresponding value or values of y can, generally, be calculated by common arithmetic. The reader will here notice the reason for saying value or values of y , as the extraction of the square

root, giving positive and negative values, will lead to two different values of y for the one value of x .

If the statement of an algebraic function is free from roots and contains only the operations of addition, subtraction, multiplication, division, and evolution, it is said to be a rational function. $y = \frac{a+x}{b-x^2}$ is such a function. Assigning a value to x , no root extraction or evolution is required for the evolution of y . Further, for each value of x there is only one value of y corresponding. A function of this kind is said to be single-valued.

A rational function composed entirely of the algebraic sum of powers of x and constants is called an integral function, $y = a + bx^2 + cx^3$, is such a function.

We were careful to state above that the value of an algebraic function of x for any assigned value of the independent variable can generally be calculated. We now give an illustration of the need for this qualification. Consider the rational function $y = \frac{4-x^2}{2-x}$. Let us find the value of y when $x=2$ in the usual way. Putting 4 for x^2 and 2 for x , the function becomes $y = \frac{0}{0}$. When an answer of this kind is obtained y is said to be undetermined, and it is evident that no numerical value can be directly assigned to y , for as any number multiplied by 0 gives 0 as an answer, $y = \frac{0}{0}$ implies that y can have any value we please. If we examine the function we see that the numerator of the fraction $4-x^2$, is the difference of two squares, and is therefore $(2+x)(2-x)$ so that the function is the same as $\frac{(2+x)(2-x)}{(2-x)}$. Let us now suppose that x has a value a little less than 2 by the small quantity h , or that $x = 2 - h$ and consequently $h = 2 - x$. The function now becomes $\frac{(4-h) \times h}{h}$, or $(4-h) \times \frac{h}{h}$. Now, however nearly x may approach 2, or the smaller h may become, $\frac{h}{h}$ must be equal to 1 as long as h has any value whatever. Thus the function has the value

$4 - h$, for all values of h , and the smaller h becomes the nearer does y approach the value 4. This fact is often expressed symbolically $\lim_{x \rightarrow 2} y = 4$, which is read, the limit of y as x approaches 2 is equal to 4. A similar argument would show that as x , decreased and approached the value 2, the limiting value of y would be 4 also.

Some readers may consider that all the foregoing is merely a roundabout way of arriving at a result obtained simply and obviously by saying that $\frac{(2+x)(2-x)}{2-x}$ must be equal to $2+x$. It must be realised, however, that to remove a common factor like $(2-x)$ by what is commonly called cancelling is mathematically legitimate only when this factor is not equal to zero, and we have pointed out as early as on p. 16 that the operation of dividing by zero is meaningless. We are only entitled to cancel the $2-x$ factor subject to the condition that it is not equal to zero. When x is exactly equal to 2 the original function is undefined, but as it continually approaches the value 4 as x approaches 2 from above or below we conclude that 4 is the value corresponding to $x=2$. In other words, we can make the defined value of y differ from 4 by as little as we please, by making x differ from 2 by the same amount.

Continuity. Let us consider one of the simplest of the rational functions $y = \frac{a}{x}$; suppose $a=1$ so that the function take the even simpler form $y = \frac{1}{x}$. This is sometimes the reciprocal function because y is the reciprocal of x . It also expresses what is sometimes called inverse proportionality. It is evident that as x increases y decreases, when $x=1$, $y=1$, and when $x=-1$, $y=-1$. If x has a very small positive value, say 10^{-8} , y has a correspondingly large value positive, 10^8 . When x has the very small negative value -10^{-8} y has the large negative value -10^8 . What is the nature of the variation of y , as x diminishes from 10^{-8} , through zero, to the value -10^{-8} . The nearer x approaches 0 from above the greater y is positive; immediately the value 0 is passed by x , y takes an exceedingly high negative value.

When x is exactly 0, $y = \frac{1}{0}$ by the algebraic formula, but this has no precise meaning, although it is commonly said that y becomes infinite. Without trying to elucidate the meaning of $\frac{1}{0}$ we may say that a change of x as small as we please from one side of zero to the other causes an exceedingly large change of y , and the smaller the change of x the greater is the change of y . y is said to be discontinuous at the point $x = 0$.

A function may be said, roughly, to be continuous as long as the change in the dependent variable y can be reduced by reducing the change in the independent variable x . A more precise definition of continuity is this: if a change of the independent variable x , by an amount h , produces a change of the dependent variable k , then k can be made smaller than any assigned quantity by reducing the change h , if the function is continuous.

To illustrate this, let us consider the function $y = x^{1000}$, where y is equal to x raised to a very high power, and where, when x is large, a small change in its value will give a much greater change in the value of y . Suppose x changes to $(x+h)$, y will change from x^{1000} to $(x+h)^{1000}$, and we have seen on p. 23 that, if this last expression be expanded the first term will be x^{1000} , the second will be $1000x^{999} \times h$, and the third term will contain h^2 and will, if h is very small, be negligible compared with the second. The change in y will, therefore, be $1000x^{999}$ times the change h in x . Now, although $1000 \times x^{999}$ may be an exceedingly large number, h can always be diminished so that $1000x^{999} \times h$ can be made as small as may be desired; hence the function is everywhere continuous. This argument shows, as a matter of fact, that all integral functions of x are everywhere continuous.

Functions used in applied mathematics are normally continuous, although there may be discontinuities for special values of the independent variable, as in the function $y = \frac{1}{x-a}$, which is discontinuous at the point $x = a$, at which the denominator of the fraction becomes zero. It is, however, possible to have functions that are everywhere discontinuous. Take,

for example, the function of x which is equal to 2 when x is a rational number, and to 3 when x is otherwise. This, it must be realised, is a true function because in assigning any value to x the value of the function is known precisely, but the function is everywhere discontinuous because everywhere an indefinitely small change in x from a rational to an irrational number will cause the value of the function to alter from 2 to 3 or from 3 to 2, and this alteration cannot be reduced.

An important corollary of the definition of continuous functions should here be noted. If $y=f(x)$ is everywhere continuous and if two values of x , say x_1 and x_2 make y respectively positive and negative, then between x_1 and x_2 there is at least one value of x that will make y equal to zero. This of course does not apply to functions that may be discontinuous for, if $y=\frac{a}{x}$ then $x=1$ makes y positive and equal

to $+a$, and $x=-1$ makes y equal to $-a$, but no value of x between $+1$ and -1 can make y zero. The foregoing can be stated in the following form: if $y=f(x)$ is everywhere continuous, y must pass through a zero value if it changes its sign; if the function is discontinuous y may change sign by passing through an infinite value.

There is another kind of discontinuity that should be noticed in an elementary study of functions. Consider $y=\sqrt{4-x^2}$; we see that if $x=0$ $y=\pm 2$, if $x=1$ $y=\pm\sqrt{3}$, and if $x=2$ y is zero. Suppose x is greater than 2, $(4-x^2)$ is negative, and y becomes the square root of a negative real number; in short the character of y changes from real to imaginary as x , increasing, passes through the value 2. If x is very slightly under 2, y will have a small value real, and if x is very slightly over 2, y will have a small value imaginary. The change in y due to a small change in x through the value 2 could, therefore, be expressed as a complex number. $k+jk$, having a very small modulus or magnitude, and this modulus could be made as small as we pleased by reducing the change of x from under to over 2. The function, however, cannot be considered to be continuous at $x=2$, because the character of the dependent variable y changes from a real to an imaginary

character. We shall see later that this has an important geometrical significance.

Transcendental Functions. A function of an independent variable that cannot be symbolically represented by a limited number of arithmetical operations is said to be transcendental. A typical example of such a function is $y = a^x$. The difference between this function and $y = x^n$ should be carefully noted; x^n means a variable number raised to a constant power, a^x means a constant number a raised to a power which may continually vary. A little thought shows that the idea of a variable power takes us into deep waters, for, as long as the value of x is rational and can be expressed in the form $\frac{m}{n}$ where m and n

are whole numbers, a^x has a precise arithmetical meaning, the n th root of the m th power of a , but if x is irrational no arithmetical meaning can be assigned to the function, for an irrational index cannot be the equivalent of extracting the root of a power. Let us consider this point more carefully. Take the function 10^x , and suppose $x = \sqrt{2}$, what meaning are we to assign to $10^{\sqrt{2}}$? Now $\sqrt{2}$ has a value between 1.414 and 1.415. $10^{1.414}$ is the same as $10^{\frac{1414}{1000}}$, it means the 1000th root of the 1414th power of 10, and as a matter of fact its value is about 26. Similarly $10^{1.415}$ means the 1000th root of the 1415th power of 10, and is a little greater than the root of the smaller power. If we could be sure that the function a^x is everywhere continuous we could say with certainty that the value of $10^{\sqrt{2}}$ lies between two rational powers of 10, the index of one being very slightly less and that of the other being very slightly greater than $\sqrt{2}$, and that, by diminishing the difference between these two rational indices, we could be sure that $10^{\sqrt{2}}$ had a value lying between two intelligible values of the function 10^x having a difference smaller than any assignable magnitude. This assumption, however, rests on the question of the continuity of the function a^x , for it is conceivable that it might be discontinuous for irrational values of x , although there appears to be no reason why this should be so. We shall see, later, how from geometrical considerations we can infer the continuity of a function of

this kind and so, assign a meaning to the expression a^x , for any value of x , rational or irrational.

We have already seen that if $y = 10^x$, x is called the common logarithm of y , and this leads to an inverse kind of function generally expressed $x = \log_{10} y$. a^x is called the exponential function, and $y = \log_{10} x$ is called a logarithm function. Assigning a value, say 2, to x , we have seen that it will be possible to find some root of some power of 10 that will approximate as closely as we please to the number 2, and, if the function is continuous, we may conclude that there is some value of y , which, used in the exponential function 10^y , will correspond to an answer exactly 2, by following an argument similar to that used in the preceding paragraph.

The trigonometrical ratios defined in Chapter III are examples of transcendental functions. $y = \sin x$ expresses a functional relationship between the arc x of a circle of unit radius and a particular side y of a right-angled triangle connected with the arc. Assigning a value to the arc x , we are sure that there is a value of y corresponding, and we could by a geometrical construction obtain an approximation to this value. But, at present, we have no idea how the y corresponding to a stipulated value of x could be worked out by an arithmetical process and without a geometrical construction.

Numerical Evaluation of Functions—Series. The final end of applied mathematics is the evaluation of numerical results, and however useful a function may be as a quantitative representation of some physical phenomenon, the knowledge embodied therein cannot be put to actual practical use unless numerical values of the function can be calculated for assigned values of the independent variable. For this reason a most important branch of mathematics is that of finding algebraic expressions, easily amenable to the elementary processes of arithmetic, which correspond in some way with mathematical functions of higher kinds than algebraic, and which cannot be exactly represented by a limited number of arithmetical operations. As an illustration of this point we may cite a particular form of the exponential function, e^x , which represents a special number e , having a value about 2.71 raised to a

variable power x . The numerical value of e^x can be represented in the following way :

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} \dots$$

and so on, continued indefinitely. The factorials, like $3! = 3 \times 2 \times 1$, have the values explained on p. 17, the powers of x increase by 1 as we proceed along the unlimited series of terms, and the index of x and the order of the factorial are the same. It is evident at once from an inspection of this equation that the operation of involving a number e , which as a matter of fact is irrational, to any power is stated as being equivalent to the addition of an unlimited number of terms each of which can be evaluated arithmetically. It is essential to realise clearly what the equation implies. The implication is this : that if we evaluate, say e^3 , by the arithmetical process represented by the right-hand side, and then similarly evaluate, say e^2 , and multiply these results together we shall obtain as an answer the number we should arrive at by proceeding directly to evaluate e^5 , for e^5 is equal to $e^3 \times e^2$.

But what do we mean by saying that the number represented by e^2 , say, which is round about $7\frac{1}{2}$ numerically, can be the equivalent of an unlimited number of fractions of the form $\frac{x^n}{n!}$? How many of these terms have we to calculate and add together to obtain a practically close approximation to the actual value of e^2 ? and what error is there in stopping in the calculation and neglecting the unlimited number of fractions which follow those we have taken into account? These are important and far-reaching questions which penetrate deeply into the theory of pure mathematics, but which we shall deal with in an approximate and simple way.

A collection of simple algebraic terms or functions like those on the right-hand side of the equation we have been considering in which each term is formed from the one preceding it according to a definite and invariable law, is called a progression or series. In the series for e^x we see by a little thought that, as the fourth term is obtained by multiplying

the third term by $\frac{x}{3}$, the term after the n th is obtained by multiplying this n th term by $\frac{x}{n}$.

Arithmetical and Geometrical Series. The simplest of all series is that generally known as an arithmetical progression, in which neighbouring terms have the same difference. An example of such a series is the progression of odd numbers 1, 3, 5, 7, and so on. We can find the sum of any number of the terms in such a series without actual addition in the following way. Denote the sum by S , and let us take six terms for simplicity, then

$$S = 1 + 3 + 5 + 7 + 9 + 11$$

and re-writing the series in the reverse order we have

$$S = 11 + 9 + 7 + 5 + 3 + 1$$

and adding these two statements together,

$$2S = 12 + 12 + 12 + 12 + 12 + 12 = 6 \times 12$$

so that $S = 6 \times 6$.

A little study of this calculation shows that, however many terms we take, say n , the last term in the direct progression is $2n - 1$, so that the sum $2S$ is equal to $n \times 2n$ and $S = n^2$. The sum of n odd numbers starting at 1 is therefore n^2 . The reader ought to be able to illustrate this geometrically by starting with one square and placing three squares round it to make a larger square, and so on.

The foregoing is of no practical importance excepting perhaps as an illustration of the kind of tricks that mathematicians have had to devise to obtain the sum of the terms of a series. We see at once that the more terms we take of an arithmetical series the greater this sum will become.

The next simplest series is the geometrical, in which each term bears the same ratio to the one that precedes it. Thus, denoting this common ratio by r , an illustration of a geometrical series is 1, r , r^2 , r^3 , and so on. The sum of any assigned number of terms of a geometrical series can be found without direct calculation and addition by means of a formula obtained by a trick or device of the same kind as that used for

the arithmetical series. Let us denote the sum by S , and consider five terms, then

$$S = 1 + r + r^2 + r^3 + r^4.$$

Multiply each side of this equation by r , and we obtain a new equation,

$$Sr = r + r^2 + r^3 + r^4 + r^5.$$

Now subtract the right-hand side of the second from that of the first, and so with the left-hand sides. In the right-hand side subtraction we see that r subtracts from r leaving zero remainder, and similarly with r^2 , r^3 , and r^4 , so that the answer is

$$S - Sr = S(1 - r) = 1 - r^5, \text{ so that } S = \frac{1 - r^5}{1 - r}$$

and it is easy to see that if we had taken any number of terms, n , the answer would have been

$$S = \frac{1 - r^n}{1 - r} \text{ which, of course, is the same as } \frac{r^n - 1}{r - 1}.$$

This means that, if r is equal to, say, 2, then the sum of six terms is, as $r - 1 = 1$, equal to $2^6 - 1$ or $64 - 1 = 63$. If $r = \frac{1}{2}$, then the denominator $1 - r$ is $\frac{1}{2}$, and the sum is $2(1 - \frac{1}{64}) = \frac{126}{64}$.

Sum to Infinity—Convergence. It is at once evident that if the common ratio r of a geometrical series is greater than 1, then, however small the first term may be, the numerical value of the terms gradually increases, and the sum can be made as large as we please by taking a sufficient number of terms. If however r is less than 1, the terms decrease steadily as we proceed along the series. The more terms we take the greater will be the sum, certainly, but can we make this sum greater than any number we like to assign? To answer this question let us consider the simple geometrical series

$$1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} \dots$$

in which each term is half the preceding. We can easily find an objective meaning for the adding up of an indefinite number of these terms. Consider a line 2 feet long, take $\frac{1}{2}$ of it, 1 foot remains, take $\frac{1}{2}$ of this remainder $\frac{1}{2}$, and $\frac{1}{2}$ remains, take $\frac{1}{2}$ of this second remainder, $\frac{1}{4}$, and $\frac{1}{4}$ remains. The summation of the series is, therefore, equivalent to progressing

along the line by steps each equal to half of what is still required to go the whole 2 feet. We have here the old classical paradox of Achilles and the tortoise in another form. However far we progress along the line by these successively diminishing steps we can never reach the end or 2 foot mark, but the more steps we take the nearer we shall approach this mark. We, therefore, say that the limiting value of the sum of the series as the number of terms is increased indefinitely, is exactly 2. We can never actually reach 2, but by taking sufficient terms we can make the difference between the actual sum and the limiting value as small as we please. A limiting value of the sum of a series of this kind is sometimes called the sum to infinity.

Now let us see how our formula fits in with the common-sense conclusion we have arrived at. As r is $\frac{1}{2}$ the sum of n terms is $2\left(1 - \frac{1}{2^n}\right)$ or $2 - \frac{2}{2^n}$. This confirms our argument, for by taking n large enough, we can make 2^n , the denominator of the fractional defect from 2, as large as we please; by taking the sum of 1,000,000 terms, the denominator of the fraction giving the difference from 2 would require about 300,000 figures to express it.

A series which has the property that the sum of its terms has a limiting value which may be approached as nearly as we please by adding a sufficient number is said to be convergent. A series in which the sum may be made to exceed any assigned value by taking sufficient terms is said to be divergent.

Let us examine the geometrical series a little more closely in regard to this matter of convergence. We have seen that when the common ratio r is greater than 1 the series is divergent and it is easy to see that this is the case when r is 1, and all the terms are of the same value. Suppose r is less than 1, then the numerator of the fraction giving the sum is $(1 - r^n)$. Now a fraction less than 1 multiplied by itself becomes smaller, however near to 1 it may be, and successive multiplications, or involution, steadily reduce the value. We conclude that, provided r is less than 1 by ever so small amount, r^n can be made as small as we please by taking n sufficiently large.

This means that the sum to infinity is $\frac{1}{1-r}$; a definite and calculable number, so that the series is convergent subject to this condition, namely, the common ratio is less than 1. The same is true if the series, instead of starting at 1 as we have taken for the sake of simplicity, starts with any other a , for, as this will be a common factor of all the terms, the sum to infinity will be simply a times $\frac{1}{1-r}$ or $\frac{a}{1-r}$.

Another example of a convergent geometrical series is

$$1 + \frac{1}{10} + \frac{1}{100} + \frac{1}{1000} \dots$$

which is a long way of writing the circulating decimal 1.111... or 1. $\dot{1}$. Here $r = \frac{1}{10}$ and the sum to infinity is 1 divided by $1 - \frac{1}{10}$ or $\frac{9}{10} = 1\frac{1}{9}$, which is the value obtained by converting the circulating decimal to a vulgar fraction by the rules of arithmetic. These rules, in their general form, are based upon the summation to infinity of a geometrical series, and the reader will find it a useful mental exercise to verify the rules, by using the general formula we have deduced, if he cares to do so.

The convergency of the geometrical series for values of r less than 1 has been established by the deduction of a formula giving the sum of any required number of terms. It is not possible to deduce such a formula for the great majority of series, and their convergency or otherwise has to be tested indirectly. The importance of this matter of convergency is manifest; for, if a transcendental function is purported to be represented by an infinite or unlimited series, and this function is known from other considerations to have defined values, the representation can only be legitimate in the practical sense if the series is convergent and has a true limiting value.

It might be thought that any series is convergent if its terms steadily and continually decrease, but this is not so. The so-called harmonic series, $1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{5}$ has continually decreasing terms, and by proceeding far enough we can find a term that is less than any assignable number, yet by adding sufficient terms we can obtain a sum greater than any assignable number, and the series is convergent. This may be

shown roughly in the following way : the sum of the 3rd and 4th terms is greater than $\frac{1}{2}$, so is the sum of the next four, the sum of the next eight and so on, so we can bracket off the series into an unending series of groups each of which is greater than $\frac{1}{2}$, so the sum can have no limit although we might have to take millions of terms to give a sum of 50.

If, however, a series with continually decreasing terms has these terms alternately positive and negative then the series is convergent. $1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} \dots$ for instance is a convergent series, and has a limiting sum, S , which is such that $e^S = 2$. It is one of the paradoxes of mathematics that the sum of this series is altered by rearranging the order of the terms. We can illustrate this quite simply. Add to each negative fraction its double and the series becomes $(1 + \frac{1}{2} + \frac{1}{3} \dots)$, now subtract twice the sum of the negative fractions or $2(\frac{1}{2} + \frac{1}{4} + \frac{1}{6} \dots)$, or $(1 + \frac{1}{2} + \frac{1}{3} \dots)$, which should cancel adding the double of each negative fraction to the original series. The sum according to this calculation is zero, and not the true sum which is a number between $\frac{1}{2}$ and 1. This, however, is a matter of mere interesting curiosity and not one of practical moment to the engineer.

When all terms of a series are either positive or negative its convergency cannot be found by mere inspection. The complete tests for convergency belong to advanced pure mathematics ; we shall content ourselves here by mentioning the simplest of these tests. It depends fundamentally on using the converging geometrical series as a kind of standard of reference, for if it can be shown that if the terms of a series are all less than the corresponding terms of a converging geometrical series, then the series can itself be inferred to be convergent. This leads to the rule that if after a fixed term the ratio of every term to the one preceding it is less than unity, then the series is convergent.

Let us apply this rule to the e^x series referred to above. We have seen that the term after the n th is obtained by multiplying this term by $\frac{x}{n}$. Now, however large x may be,

there can always be found the term for which n is just greater than x , and after this term the ratio of every succeeding term

to the one preceding is always less than 1, so that the series is convergent for all values of x . This is in some ways a remarkable conclusion, and a little time can well be spent in pondering on it.

The utility of a convergent series for practical calculation depends a good deal on the rapidity of convergence, that is, upon the number of terms that have to be taken to obtain a sufficient approximation to the limiting value for practical purposes, and much ingenuity has been exercised in the devising of rapidly converging series for the numerical calculation of functions that are frequently wanted in applied science.

To illustrate the convergence of a series we may give the arithmetical calculation of the value of e^2 from the series

$$1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} \dots$$

by putting 2 for x : The calculation is:

1· 000 000	× 2
2· 000 000	× 2 ÷ 2
2· 000 000	× 2 ÷ 3
1· 333 333	× 2 ÷ 4
666 667	× 2 ÷ 5
266 666	× 2 ÷ 6
088 888	× 2 ÷ 7
25 397	× 2 ÷ 8
6 349	× 2 ÷ 9
1 432	× 2 ÷ 10
284	× 2 ÷ 11
52	× 2 ÷ 12
8	× 2 ÷ 13
1	
<hr/>	
7· 389 063	

of which the last figure certainly will be doubtful. It will be seen that after the 10th term the series begins to converge rapidly. Convergence commences after the third term.

In this chapter we have been concerned merely with a brief

and elementary sketch of the idea of function, considered from the algebraic or symbolic standpoint, and we shall have to leave till a later stage a detailed study of the simpler transcendental functions used in engineering. In the following chapter we shall study functional relationship in a different way, and we shall see that the method of graphical representation sheds much light on what is not very clearly revealed by algebraic representation.

CHAPTER VI

GRAPHS

Basic Considerations. In Chapter IV we studied at some length two methods of symbolically representing a directed step from a reference point on an unlimited straight line, and we saw that this step, or what is the same thing, the end of the line showing it, can be defined in two ways: the first by the distance of the end point P from the reference point, O (Fig. 27), and the angle the line OP makes with the unlimited horizontal reference line; and, secondly, by the length of

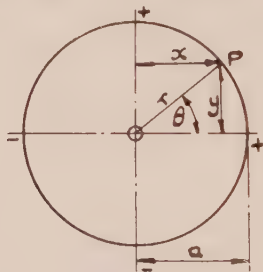


FIG. 27.

two directed steps the one horizontal x , and the other vertical y which combined will make up the whole step represented by the line OP . If we are considering merely the physical position of the point P , we can say that this position is defined by the perpendicular distances from two unlimited lines at right-angles to each other and intersecting at the point O ; the perpendicular distances from the vertical line, x , are positive to the right and negative to the left of it, and the perpendicular distances from the horizontal line, y , are positive above and negative below it. Having fixed the two reference lines intersecting at O , the position of a point P can be completely specified in these two ways, and by two sets of magnitudes, the one r and θ , and the other x and y . These sets of magnitudes are called respectively the co-ordinates of the

point P ; r and θ are the polar co-ordinates in reference to a fixed line, x and y are the cartesian co-ordinates in reference to two lines at right-angles called the axis of co-ordinates. The x or horizontal distance defining the position of P is sometimes called its abscissa, and the y or vertical distance its ordinate. The reader should note that we are now thinking of the position of a point, and not of a directed step, and x and y are actual geometrical lengths and not the objective representations of two different kinds of numbers.

Suppose we draw a circle of radius a with the reference point O as centre, the circumference of the circle is an example of what in mathematics is called a locus ; each point of it satisfies a definite condition. What is this condition for the circle ? Evidently that each point is the same distance from the reference point O. We express this symbolically by $r=a$, and we must consider carefully what this statement, as a piece of mathematical symbolism, implies. The position of a point is generally fixed by two magnitudes, θ as well as r , but the statement $r=a$ says nothing about θ . This means simply that θ is irrelevant ; it has nothing to do with r , which is the same, whatever value θ may have, which, considered from the geometrical point of view, is evidently true.

We can, however, express the condition satisfied by the circular locus in another way, for the radius, the x , and the y of any point define a right-angled triangle of which the hypotenuse is equal to a constant length a . We express this symbolically by $x^2 + y^2 = a^2$, and we observe that, whether x and y are positive or negative, their squares are always positive, so that they give a positive value to a^2 . But $x^2 + y^2 = a^2$ is evidently like an algebraic function connecting two variables x and y , and we can easily change this statement and obtain an explicit function of y by writing, first, $y^2 = a^2 - x^2$, and, finally, $y = \pm \sqrt{a^2 - x^2}$. What does this mean ? Primarily, that taking any point on the curve for which the value of the x co-ordinate is known, and finding the value of y by arithmetical calculation, the answer will be the actual y ordinate of the point. But the whole meaning is a little more subtle than this, for we have to account for the significance of the double sign \pm , which has been used because the square

root of a positive number has two equal values, the one positive and the other negative. A fuller geometrical meaning of the function $y = \pm \sqrt{(a^2 - x^2)}$ is, that giving x any value positive or negative, numerically less than a , and evaluating the two equal values of y by means of the function, the chosen value of x with the equal positive and negative values of y will locate two points on the circumference of the circle. We make the qualification that x must not be greater than a , because, otherwise $a^2 - x^2$ will be negative, its square root will not be a real number, and the geometrical vertical distances represented by the y 's of the circle must be expressed by real numbers.

Let us state explicitly two mathematical relations between the function $y = \pm \sqrt{(a^2 - x^2)}$ and a circle of radius a with the reference point O as centre. First, having defined the circle geometrically as a locus, we have found that the cartesian co-ordinates of all of its points are connected by the functional relationship $y = \pm \sqrt{(a^2 - x^2)}$ so that, selecting any point, and finding the y by calculation from the formula or function from the x value of this point, the answer to the sum gives the actual y value. From this point of view the function $y = \pm \sqrt{(a^2 - x^2)}$ is said to be the equation of the curve, or the circumference of the circle, as it expresses algebraically the geometrical criterion by which the curve is constructed. Secondly, taking the function $y = \pm \sqrt{(a^2 - x^2)}$ and giving x all values positive and negative less than a , and calculating the values of y corresponding, we shall obtain a collection of sets of y and x values each of which set will be a point of the circumference of the circle. From this point of view the circle is called a graph of the function, and locating points of the graph by first calculating sets of corresponding x and y values is called plotting the graph. Everybody knows how the process of plotting is facilitated by the use of squared or graph paper.

The reader should realise that it is no mere tautology to distinguish, as we have done, between these mutually inverse relations between a mathematical function and the positions of the points on a geometrical line. In applied science we sometimes require not only to express as a function, or

analytically as it is often termed, the characteristics of a locus precisely defined in the geometrical sense, but we often arrive at graphs which represent the results of experiment or observation for which a functional relationship in the analytical sense has to be found. Sometimes even the curve or graph is described by the instrument used in the experiment; an oscillograph, for instance, records a graph. Another important part of applied mathematics is that of plotting or representing graphically an exact functional relationship deduced analytically by theory, we order, as we shall shortly explain, that the characteristics of the function may be concisely and pictorially represented. Thus, in our basic illustration, an actual geometrical circle shows at once what is only implicitly stated in the function $y = \pm \sqrt{(a^2 - x^2)}$, that the maximum values of x and y are each a , and that the points represented by sets of corresponding values are equidistant from the reference or origin point O. We shall commence our study with this branch of mathematics from the second of these two stand-points, that of the graphs of defined analytical functions.

Graphs of Linear Functions. Let us consider the simplest of all functions connecting the y with the x of a graph, that of direct proportionality and expressed as $y = nx$. We see at the outset that as x and y both represent lengths of lines on a diagram, n must be a ratio or pure number, and that the function is really $\frac{y}{x} = n$ a constant. Further, as when $x = 0$, $y = 0$, the graph must pass through the reference point or origin. As the x and y of any point determine a right-angled triangle and as the ratio $\frac{y}{x}$ is constant, all these triangles must be of the same shape, so that the graph is like the straight line L in Fig. 28, extending indefinitely in both directions, the downward direction from O, it will be seen, giving negative values to both x and y , so that their ratio remains positive and equal to n . n is the rise of the graph, going to the right, for a unit increase of the horizontal distance x ; it is called the slope, and is reckoned positive when, as with the line L, the ratio $\frac{y}{x}$ is positive. The line L₂ in Fig. 28 has a negative

slope and would have an equation $y = -mx$. Thus proceeding along a graph to the right from the origin point O, the slope is positive for a rise and negative for a fall. How can we define this slope geometrically. Since $\frac{y}{x}$ is the tangent of the angle the line makes with the horizontal reference line or x axis, and is constant and equal to n , then as $\tan \theta = n$ we can say that $\theta = \arctan n$. If the reader is not familiar with this notation he should turn back to p. 50. $\theta = \arctan n$ is what is called the polar equation of the straight line L,

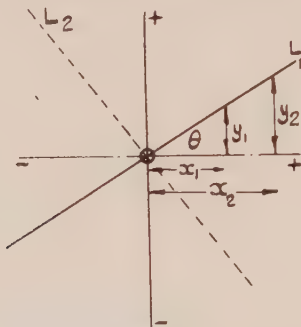


FIG. 28.

because it expresses the relation of the distance r of any point on it to the inclination of the line from O to this point. But r does not appear in the equation. This means that it is irrelevant, and that for any value of r at any point in the line θ is constant, which is evidently true, for the distance r is measured along the line.

Now let us consider the slightly more complicated function $y = b + nx$. We observe that the quantity b must be a definite length on the graph, and that, as a matter of fact it locates the point where the value of x is zero. We see also that $y - b$ is directly proportional to x . The graph of the function is shown in Fig. 29; it is a straight line cutting the vertical y axis at a point distant b from O, and with a slope equal to n . It is evident at once from an inspection of the diagram that $y - b$ is proportional to x , for at any point, $y - b$

and x determine right-angled triangles, all of which have the same shape.

We can alter the equation $y = b + nx$, for the ratio n is equal to the ratio of b to the reverse of the distance a on the x axis

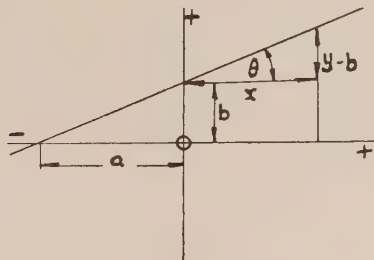


FIG. 29.

from O , where the graph intersects this axis. The slope is not equal to $\frac{b}{a}$ for a is negative and the slope is upwards or positive; it is equal to $-\frac{b}{a}$. The equation is, therefore,

$y = b - \frac{bx}{a}$ or dividing all terms by b , and transposing the x

term $\frac{x}{a} + \frac{y}{b} = 1$. This is called the intercept form of the

equation because the distances a and b from O where the graph crosses the axis are known as the intercepts. This latter form is important in pure mathematics, but in applied science the equation to a graph like that of Fig. 29 would be expressed as $y = b + nx$ in which y is an explicit function of y . A little thought will show that we can have four kinds of graphs of the kind shown in Fig. 29. b and n can be both positive or both negative, and b can be positive and n negative and vice versa. The reader should think this out and draw the three types of other than that shown in Fig. 29 for which both b and n are positive. If he does this he will see that these four types comprise all possible kinds of straight line graphs. We conclude, therefore, that all such graphs can be represented by an explicit function like $y = b + nx$, where b and n may be either

positive or negative. Such an explicit function, in which only the first power of the independent variable x appears is therefore often called a linear function. The reader should be careful to distinguish between a linear function of y and a particular case of this function, direct proportionality, when $b=0$ and $y=nx$.

A little thought will make it quite clear that the inclination θ of a line joining a point on the graph of Fig. 29 to the origin O is continually changing as the length r of this line alters. We cannot, therefore, find a polar equation to the graph so simple as that deduced for Fig. 28. As was shown on p. 53, the x of the Fig. 29 graph is equal to $r \cos \theta$ and the y to $r \sin \theta$ so that the equation $y=b+nx$ can be changed to $r \sin \theta = b + rn \cos \theta$, which can be transformed to give r as an explicit function of θ . But this function, containing the transcendental functions $\sin \theta$ and $\cos \theta$ is of little practical use.

This point illustrates a general kind of rule relating to graphs; that a simple kind of equation relating to x and y or cartesian co-ordinates usually leads to a complicated equation relating to r and θ , or polar co-ordinates, and vice versa. As an example of the vice versa case we may cite the very simple function $r=a\theta$, stating that r is directly proportional to θ . This is the equation to the spiral curve of Fig. 30, the cartesian

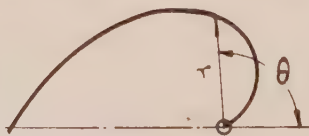


FIG. 30.

or x and y equation of which is a very complicated function. We have mentioned this point and given this simple illustration because, as equations in polar co-ordinates are not very important in the simpler applications of mathematics to engineering, we shall say no more about them, and shall confine our attention to the graphical interpretation of functions of y in terms of x .

Graphs of Simple Algebraic Functions. The next simplest function of x to consider is $y = \frac{x^2}{a}$, and here we see at once that

y and x are zero together, so that the graph passes through the origin point O . As y and x both stand for lengths a must be a length also. It is easy to see that, giving a the value 1, we get by calculation corresponding y and x values like 1 and ± 1 , 4 and ± 2 , 9 and ± 3 , and so on, so that the graph is the full line curve of Fig. 31 symmetrical about the vertical or y axis, because equal positive and negative values of x give the same value y . The meaning of the length a is also

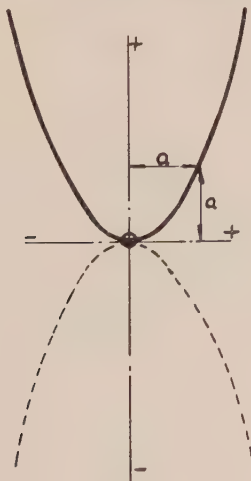


FIG. 31.

clear; it is the value of x that gives an equal numerical value to y . If a is negative, so that the equation is $y = -\frac{x^2}{a}$, the graph is the dotted curve of Fig. 31. Both of these curves are called common parabolas in geometry.

We see from the two graphs that the value of y corresponding to a given value of x is not changed if x is changed to $-x$. A function with this property is called an even function. If changing the sign of x merely changes the sign of y and does not alter its numerical value the function is called an odd function; $y = nx$ is an odd function. A function like $y = b + nx$, in which a change of sign of x alters the numerical value of y

is neither even nor odd. The graph of an even function is symmetrical about the vertical y axis. The graph of an odd function has a kind of symmetry about a line drawn through O and bisecting the angle between the axes.

The most general form of an integral function of x containing powers of x up to the second is one like $y = b + nx - \frac{x^2}{a}$. We can easily work out what the shape of the graph of this function will be by studying Fig. 32. As b is the value of y , for $x = 0$,

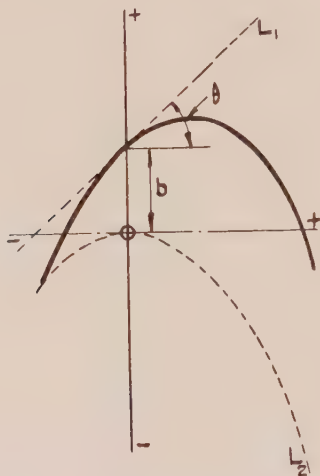


FIG. 32.

the graph cuts the y axis at a distance b from O . As a matter of fact the y of the graph is the sum in the algebraic sense of the y 's of two subsidiary graphs, the one L_1 in the figure, $y = b + nx$, and the other, L_2 , $y = -\frac{x^2}{a}$. The influence of the $y = b + nx$ subsidiary graph causes the principal graph first to rise, and that of the other graph $y = -\frac{x^2}{a}$ causes it ultimately to go downward. It can be shown that the principal graph of Fig. 32 is like the L_2 graph, a common parabola with its central line and its highest point displaced.

An important class of functions is that represented by the formula $y = x^n$. The graphs of these functions for $n = 1, 2, 3$, and 4 are shown in Fig. 33. We have already studied the cases for $n = 1$ and 2. The graph for $n = 3$ is shown, and this graph really represents the function $y = \frac{x^3}{b^2}$, for y being a length, x^3 , the cube of a length should in the geometrical sense be divided by the square of another length b . If, however, we

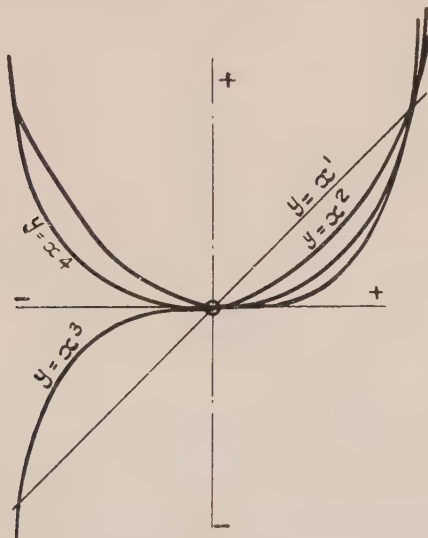


FIG. 33.

consider $b^2 = 1$, then we can deal with the function in the simple form $y = x^3$. We notice two points: first, that as with the functions $x = y^1$, and $x = y^2$, y is 1 when x is 1, secondly, that beyond this point the curve is steeper than that representing $y = x^2$, for the cube of 2 is 8, whereas the square of 2 is 4, and thirdly, that negative values of x give negative values of y . $y = x^3$, in short, is an odd function, and it is easy to see that if the exponent n of $y = x^n$ is odd, the function is odd, and, if this index is even, the function is even. The graph of the even function $y = x^4$ is also given in Fig. 33. This graph,

after the $x = 1, y = 1$ point, is passed is even steeper than that of $y = x^3$, for if $x = 2$ then x^4 is 16. The reader ought to be able to think out for himself why, after passing the (1, 1) point, the graphs get steeper as the index n is increased; the reverse is the case starting from zero. This can be understood quite well by finding the four y values, for $x = \frac{1}{2}$; these decrease as the index increases from 1 to 4. Curves of the class just considered are called parabolas.

A most important function of the $y = x^n$ class is that for which $n = -1$, which is $y = \frac{a^2}{x}$ and which represents inverse proportionality. We have already seen some peculiarities of this function on p. 93. Considering the case when $a = 1$ and the function is simply represented as $y = \frac{1}{x}$, we see that a large value of x , positive or negative, gives a small value of y of the same sign, and when x is small y is large. The function is odd, as could be inferred at once by the power of x , -1 , being odd, and the graph is as shown in Fig. 34. We see that the graph consists of two distinct branches; considering the upper one, we see also that as x decreases the graph continually gets nearer the vertical or y axis, but can never touch it, for then x being 0, a meaningless value of y would correspond. Similarly no possible value of x can make y so small that the graph, extended ever so far, can touch the horizontal or x axis.

The axes are straight lines continually approached by the graph but never actually touched or intersected by it. Such straight lines standing in this relation to a graph are called asymptotes. The graph of Fig. 34 is a rectangular hyperbola, so called because its asymptotes are at right-angles to each other. The graph is symmetrical about the line, shown dotted, bisecting the angle between the asymptotes. This graph gives a vivid pictorial representation of the anomaly or discontinuity of the function, mentioned on p. 93, that y can pass from a positive to a negative value without becoming zero.

If the reader has grasped the argument in the first part of this chapter he will see at once that the graph of the function

$y = \sqrt{2 - x^2}$ is a circle of radius $\sqrt{2}$. No real values of y can be obtained when x exceeds $\sqrt{2}$, for then $\sqrt{2 - x^2}$ becomes negative. $(2 - x^2)$ is the same as $-1 \times (x^2 - 2)$, so that the function can be written $y = \sqrt{-1 \times (x^2 - 2)}$ or $y = j\sqrt{x^2 - 2}$. The interpretation of this is that when x^2 is greater than 2 in the original function the y values are numerically equal to $\sqrt{x^2 - 2}$ but must be reckoned in a direction at right-angles to those obtained when x^2 is less than

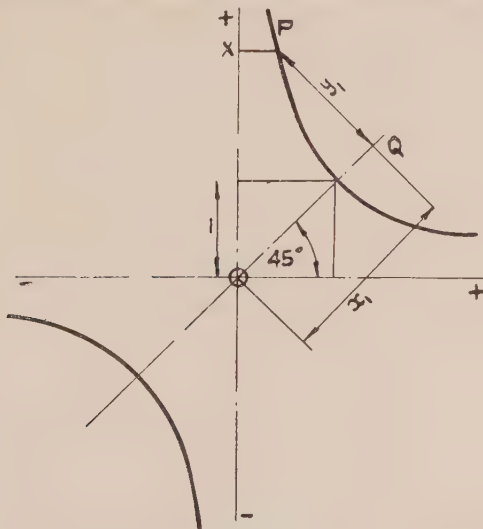


FIG. 34.

2. $y = \sqrt{x^2 - 2}$ is therefore a function represented by a graph which is a kind of extension of the circle $y = \sqrt{2 - x^2}$ in a plane perpendicular to it. Let us plot y values for x^2 greater than 2 in this plane. We obtain a graph like Fig. 35, having two similar branches which cut the x axis at points distant $\sqrt{2}$ each side of O , and which appear to be asymptotic to two lines through O , bisecting the angle between the x and y axes. We can see that these lines are asymptotic to the graph, for, the larger x is, the smaller is 2 in relation to it, and the nearer the function approximates to $y = \sqrt{x^2}$ or $y = \pm x$, which is the equation to the two lines $y = x$ and $y = -x$. The graph of Fig. 35 looks like that of Fig. 34 turned through 45 degrees,

about O, clockwise. A little concentration on the following reasoning will show that this is the case. Take a point, P, on the graph of Fig. 34, of which OY is the y , and XP the x co-ordinates. If the 45 degree line OB is to become the new x axis by rotating the graph to the Fig. 35 position, then drawing PQ at right-angles to OB, OQ is the new x co-ordinate which we call x_1 and PQ is the new y , or y_1 co-ordinate. Let us proceed along the line OQ from O and then along QP to P; we have ascended vertically a distance $OY = y$, and through a distance $x_1 + y_1$ at 45 degrees to the vertical. Fig. 12 shows

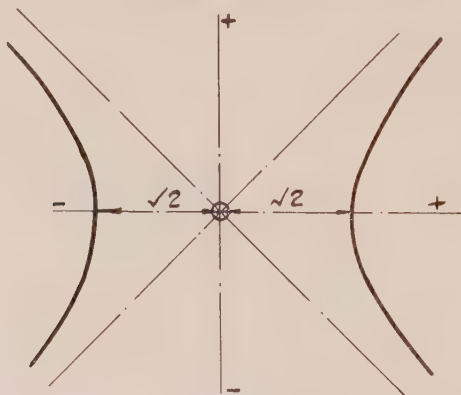


FIG. 35.

clearly, if explanation is required, that the distance $(x_1 + y_1)$ is $\sqrt{2}$ times y . Similarly in taking the OQP path we have gone, $YP = x$ horizontally to the right of O by a distance $(x_1 - y_1)$ at an angle of 45 degrees to the horizontal, $x_1 - y_1$ because going to Q has taken us too far. $(x_1 - y_1)$ is thus also equal to $\sqrt{2}$ times y . We have thus found that $(x_1 + y_1) = \sqrt{2}y$, and $(x_1 - y_1) = \sqrt{2}x$, so that $(x_1 + y_1) \times (x_1 - y_1) = x_1^2 - y_1^2 = \sqrt{2}y \times \sqrt{2}x = 2xy$. But as $y = \frac{1}{x}$, $xy = 1$. So $x_1^2 - y_1^2 = 2$, and $y_1^2 = x_1^2 - 2$, showing that the equation to the curve of Fig. 34 when twisted through 45 degrees clockwise is the function $y = \sqrt{(x^2 - 2)}$. We shall return again to this important subject of the rectangular hyperbola later when we come to study exponential and logarithmic functions.

As a final example of the graph of a simple algebraic function, let us consider $y = \frac{ax}{b+x}$. We can, without calculation, see at once the following features of the graph. When x is 0 y is 0, so that the graph passes through the origin point. As x increases, positively, the b in the denominator becomes of less and less account, and the function approximates more and more closely to $y = \frac{ax}{x}$ or $y = a$. This last equation means a line for which every point has the same y , equal to a ; it is

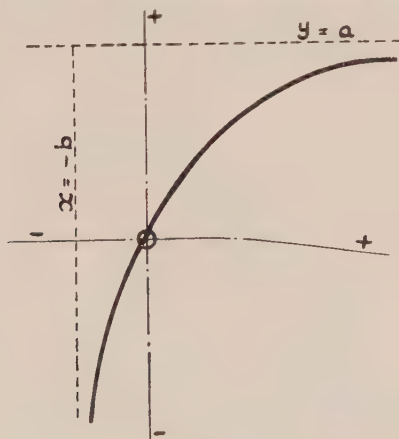


FIG. 36.

a line parallel to the x axis at a distance a above it. When x is negative and approaches the value $-b$, the denominator of the fraction will become exceedingly large; the graph is, therefore, as shown in Fig. 36; it is a rectangular hyperbola, one branch of which only is shown, with its centre shifted from the origin point O, and having for asymptotes the two lines $y = a$ and $x = -b$.

Slope Function—Tangent. A graph is a pictorial representation of the change of a dependent in relation to an independent variable that is expressed in the symbolic or functional form. All that the graph shows is implicitly contained in the function although it cannot immediately be perceived.

The outstanding characteristic of a graph is its geometrical shape which is evidently the geometrical representation of the nature of the function to which it corresponds, and the shape of a graph depends very largely upon the changes in its slope. The principal characteristic of all straight-line graphs, for instance, however they may be situated relative to the axis of co-ordinates, is that their slope is everywhere constant. A glance at the graph of Fig. 31 shows us in a general way that proceeding to the right from the origin O the slope is continually increasing, but it is not clear what meaning we are to assign to the idea of slope as applied to a curved line. This point we shall endeavour to elucidate.

Consider a point P on the graph of Fig. 37. Proceeding

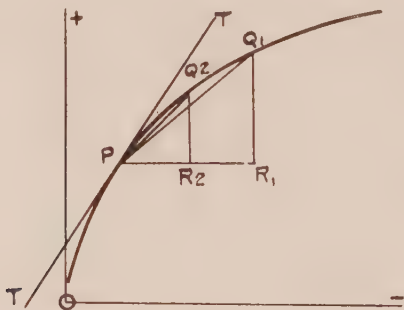


FIG. 37.

from P to the point Q_1 , we can say with certainty that the average slope in this interval is the ratio of the vertical rise R_1Q_1 to the horizontal distance traversed. This, however, is the slope of the straight line PQ_1 and not that of the curve. If we consider a smaller distance along the curve, say Q_2 , the average slope in the interval is greater, it is $\frac{R_2Q_2}{PR_2}$. Generally,

as the moving point Q approaches P , the average slope will gradually change, and it will continually approach what we have called a limiting value, that is, one that can be approached as closely as we please by making the distance PQ along the curve sufficiently small. This limiting value of average slope

is called the slope of the curve at the point P, and a straight line through P having a slope equal to that of the curve, like T in Fig. 36, is called a tangent to the curve.

Let us consider this matter algebraically in reference to the function $y = \frac{x^2}{a}$, the equation to a common parabola. For

the sake of simplicity we shall take $a=1$ so that the function is $y=x^2$. Let us think of a particular pair of values of x and y , representing a point on the parabola. If we move from this point a little away from the origin point both x and y will increase. Using the notation explained on p. 60 x will become $(x+\delta x)$ and y will become $(y+\delta y)$, but as we are still on the curve $(y+\delta y)$ must be equal to $(x+\delta x)^2$ or to $x^2+2x\delta x+(\delta x)^2$, which is the same as $y+2x\delta x+(\delta x)^2$. This means that moving from the point, along the curve, so as to increase the x distance by δx , we have ascended vertically by an amount $2x\delta x+(\delta x)^2$. The average slope of the movement is therefore the rise divided by the horizontal movement, or $\frac{2x\delta x+(\delta x)^2}{\delta x}=2x+\delta x$. We can now see, in the case of the

parabolic curve, that the slope has a limiting value at the point considered, for the average slope of a small portion of the curve beyond this point becomes more nearly equal to $2x$ as the length of the small portion, and hence the δx corresponding to it, becomes smaller and smaller, because we can make the average slope as near to $2x$ as we please by making the portion of the curve sufficiently small. We therefore say that the slope of the parabolic curve $y=x^2$ at a point having co-ordinates x and y is equal to $2x$. The reader should carefully note that we do not say that $2x$ is the slope when δx is zero, for then the average slope would be nothing divided by nothing, which gives no definite answer. We say that the average slope along a portion of the curve from the point can be made as near to $2x$ as we please by making this portion sufficiently small. The statement, slope $=2x$, for the parabola $y=x^2$ shows us that the slope of the curve is a function of x ; actually that the slope is proportional to x . We have thus obtained and deduced from the function $y=x^2$ a secondary kind of slope function which is often called a derived function.

We shall have a good deal more to say later about these derived or slope functions.

One point should be noted before we leave this discussion for the time being. The equation, slope $= 2x$ is not right geometrically, for a slope is a ratio or pure number, and x is the length of a line. Had we worked with the complete equation of the parabola, $y = \frac{x^2}{a}$, we should, as the reader can probably easily deduce for himself, have arrived at the conclusion slope $= \frac{2x}{a}$, which is correct geometrically as it gives the slope as a ratio of two lines.

Area Function. Another important kind of secondary func-

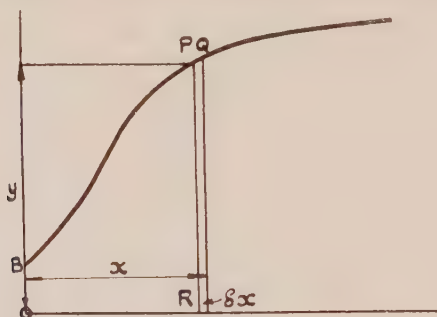


FIG. 38.

tion connected with a graph is the area enclosed by the graph, the axes of the co-ordinates, and the x axis of any point on the graph. Thus, considering the point P of the graph of Fig. 38 the area corresponding to the point P is that represented by the irregular figure $OBPR$. This area is evidently continually changing as the position of the point P changes. Now P is located by its x and y values. Consider another point Q very near to P . A shift from P to Q increases the area by a small strip, and we can write a value for this increase; it is equal to the increase of x passing from P to Q , or δx , multiplied by the average height or ordinate in the interval P to Q . If we call this increase of area δA , then by making the distance from P to Q sufficiently small we can make δA

as near to $\delta x \times y$ as we please, where y is not an average ordinate or height, but the ordinate of P. The rate at which area is increasing at the point P can therefore be said to be equal to y . Now as area is continually changing as the position of P, and the x co-ordinate area is a function of x , and it is called the area function of the graph. It would evidently be possible to plot a secondary graph of this secondary area function, and the slope of this secondary graph, at any point, being the limiting value of ratio change of area to change of x would be represented by the y of the principal or primary graph at that point. Otherwise a graph is a slope curve of the graph of its area function.

This point is illustrated in Fig. 39. The principal function

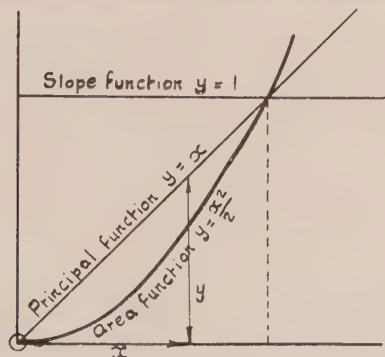


FIG. 39.

is $y = x$ represented by a straight line through the origin O at 45 degrees to the x axis. The x and y of this graph are always equal, and its slope is equal to 1, for going along the graph the rise is always equal to the horizontal travel. The slope curve is, therefore, the horizontal line $y = 1$. We can easily discover the area function, for at any point this is the area of a right-angled triangle which is half the square on x or $\frac{1}{2}x^2$. The area function is therefore $y_1 = \frac{1}{2}x^2$, where y_1 denotes area, and the graph of this function is the parabola shown. We have calculated independently that the slope of $y = x^2$ at a point determined by x is $2x$ and we can infer that when y is equal to $\frac{1}{2}x^2$ the slope will be halved and

equal to x so that the slope curve of $y = \frac{1}{2}x^2$ is given by the principal graph of Fig. 38 in which the y is always equal to the corresponding x .

Another example of the direct calculation of an area function is that of a circle which is the graph of the function $y = \sqrt{(a^2 - x^2)}$. This area for an assigned value of x is the figure shaded in Fig. 40, and it is made up of the sector containing the angle ϕ , and the right-angled triangle of height

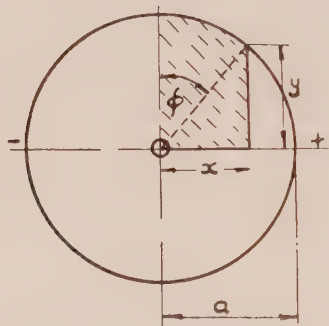


FIG. 40.

y and base x . The area of the triangle is $\frac{1}{2}xy$ or $\frac{1}{2}x\sqrt{(a^2 - x^2)}$. What is the area of the sector? If the angle ϕ is in circular measure, then the arc is ϕ times the radius a , or ϕa so that this sectorial area, as we have seen on p. 48, is $\frac{1}{2}a^2\phi$. But we see that the sine of ϕ is equal to the ratio of x to the radius a , or $\frac{x}{a}$, so that we indicate ϕ as $\text{arc sin } \frac{x}{a}$. The area function of a circle represented by $y = \sqrt{(a^2 - x^2)}$, determined by an assigned value of x is therefore $\frac{1}{2}x\sqrt{(a^2 - x^2)} + \frac{1}{2}a^2 \text{ arc sin } \frac{x}{a}$.

This looks very complicated, and had it been deduced algebraically its significance would not have been clear. Actually it stands for the sum of two component areas, the one a right-angled triangle of side x and hypotenuse a , and the other a sectorial area of radius a and with an angle $\text{arc sin } \frac{x}{a}$.

determined by the shape of the triangle. This is an example of what often applies in mathematics, that symbolic statements are not always as formidable as they appear at first sight.

Maximum and Minimum Values—Inflecting Points. Consider the graph of Fig. 41. Proceeding along it from the vertical, y axis, the slope is positive or upward for a certain interval ending at the point A at which no further increase of y takes place. The slope here is zero and the tangent is parallel to the x axis, or horizontal. The value of y at the point where the slope is zero, that is where y ceases to increase and begins to decrease is called a maximum value, not *the*

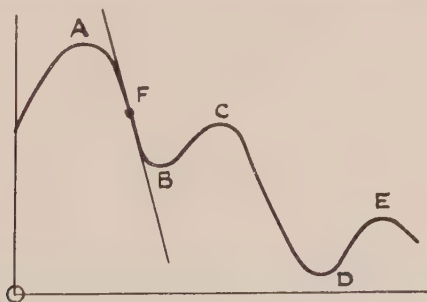


FIG. 41.

maximum value, but *a* maximum, because the y may exceed this value at some other point of the graph. When proceeding along the graph y is diminishing and the slope is negative, then, if we arrive at a point like B, where diminution of y ceases and increase starts, the value of y is called a minimum. Note, by comparing the minimum point B with the maximum point E that a minimum value may be greater than a maximum value. The terms maximum and minimum refer to values of y in the neighbourhood of points where the slope of the graph is zero. In the neighbourhood of a maximum value the slope is decreasing, passing through zero at the maximum point and then becoming negative; in the neighbourhood of a minimum point the slope increases with x , changing through zero from negative to positive at the minimum point.

In the portion A to B of the graph the slope starts and ends at zero, and its average value is negative because the graph is going downward. At one point the slope has a maximum numerical value. This point, F in the figure, is called a point of inflection, and, as the numerical value of the slope on each side of this point is less than that of the curve and hence, of the tangent at this point, it is clear that the tangent here crosses the graph.

It is evident from the foregoing definitions that, if we could find the values of x for which the slope of the graph of a function is zero, this value of x would define either a minimum or a maximum value. For example, the graph of Fig. 32 shows that there is a maximum value. The function is $y = b + nx - \frac{x^2}{a}$. We have seen that the y of the graph is the difference of the y 's of two subordinate graphs, the one $y = b + nx$ of which the slope is n , and the other $y = \frac{x^2}{a}$ of which we have seen that the slope is $\frac{x}{a}$. We can infer that the slope of the principal graph is zero when the slopes of these two subordinate graphs are equal, and this occurs when $n = \frac{x}{a}$ or when $x = 2an$. This is the value of x that gives the maximum value of y , and putting this x value in the functional equation we find that $y = b + 2n^2a - 4n^2a = b - 2n^2a$.

Continuity of Functions. The graphical representation of functions gives a very clear idea of the otherwise subtle idea of continuity dealt with on p. 93. On recollecting what was there written it will be clear to the reader that a function is continuous where the slope of its graph is a finite quantity. The hyperbolic graph of Fig. 34 has no finite slope where $x = 0$ and is there discontinuous, but there is here a distinct break in the graph, and this is not a sufficient criterion of discontinuity. The function whose graph is shown in Fig. 42 is, in a way, discontinuous at the point x , although there is no break there, for the portion AB is vertical, the slope suddenly becomes greater than any finite number, and no

reduction of an increment on x can make the change in the y value less than the quantity represented by the straight vertical portion AB.

Plotting Graphs of Mathematical Functions. If a function of two variables is expressed mathematically in the explicit form $y=f(x)$, the graph of the function for numerical values of the constants can be drawn completely by calculating arithmetically a sufficient number of corresponding sets of x and y values and locating the points corresponding to these sets on squared graph paper. Actually this method

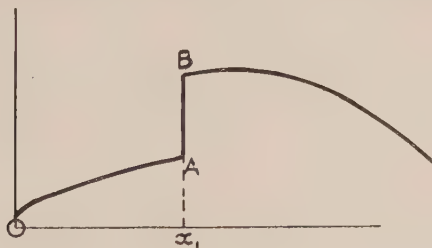


FIG. 42.

of plotting, without preliminary and subsidiary calculations, would be not only very laborious, but would be liable to considerable errors. The delineation of the graph of a function is always started by a preliminary examination to determine salient points and characteristics. Among these are the maximum and minimum points, the asymptotic branch of the graph if any, the position of these asymptotes, and any limiting values of the one variable beyond which the other becomes unreal. Thus, when we considered the function

$$y = \frac{ax}{b+x}$$

we saw that there were asymptotic branches of

the curve which continually approached the straight lines $y=a$ and $x=-b$, and, before the calculation of x and y values and the plotting of the points corresponding it would, in drawing the graph of this function, evidently be desirable first to locate and draw these asymptotes. Again, it is clear that accurate drawing of the graph will be assisted if we could find the slope and, hence, draw the tangent line at the origin

point where x and y are each zero. With the asymptotes, and this tangent, a few plotted points from corresponding calculated values will enable us to stretch the graph with considerable accuracy.

Again, consider the function $\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$, we see by a slight change of this to $\frac{y^2}{b^2} = 1 - \frac{x^2}{a^2}$, that no real values of y can be obtained when x is greater numerically than $\pm a$. Similarly from another equivalent form, $\frac{x^2}{a^2} = 1 - \frac{y^2}{b^2}$, we see that when y is numerically greater than b no real values of x are possible. The graph is a closed curve of horizontal extent from $+a$ to $-a$, and vertical extent $+b$ to $-b$. By a further change of the function we see that $\frac{a^2 y^2}{b^2} = a^2 - x^2$, which shows that if every y value is changed in the ratio $\frac{a}{b}$, the curve will become a circle, which shows that the form of the curve is elliptical. Finally, let us think of tracing the graph of a function like $y = 4 + 2x - x^2$. We know that the general shape of this graph will be like Fig. 32, and it is evidently almost essential, in delineating the graph, first to locate the maximum point and to draw the horizontal tangent there. We shall see later how maximum and minimum points can be found.

The tracing of the graphs of mathematical functions is a branch of the art or technique of mathematics which requires not only knowledge, but skill, practice, and experience, and it is so wide in its scope that one of the classics of mathematical literature ("Curve Tracing," by P. Frost) is devoted entirely to it. Although this technique is very interesting and, indeed, fascinating to those with a taste for pure mathematics, it is not much wanted by the engineer, and no more need be said about it than the slight descriptive sketch which we have given.

Empirical Graphs. When two physical quantities vary simultaneously, and sets of simultaneous values are obtained by observation and measurement, then if points are plotted on

graph paper to represent these sets and a curve is drawn to contain all the points, this curve is called an empirical graph. Primarily an empirical graph is a method of concisely combining and exhibiting the information contained in the observed data ; it does not necessarily indicate what may be called a functional relationship between the varying quantities. For instance, barometer height is continually varying, and if this height is obtained at various times and corresponding values plotted, we could hardly say that this graph would represent atmospheric pressure as a function of time, because we know that the primary cause of this pressure variation is not mere efflux of time. Many cases arise in physical science in which a true functional relationship is known to exist between two variable quantities, but a graph exhibiting a simultaneous variation obtained by experiment will not necessarily show this functional relation unless other causes of variation are eliminated. Thus, for instance, Ohm's law in electrical science states that the current in a circuit maintained at constant temperature is exactly proportional to the voltage applied to it, but a set of observed values of current and voltage might fail to plot into anything like the straight line graph implied by the law, if the circuit was of a particular type and its temperature was allowed to vary by the heating effect of the current.

If a number of points are marked on squared paper, then, theoretically, it is always possible to find a mathematical function the graph of which passes exactly through each and all of the points, but this function will not necessarily represent the physical connection of the variable magnitudes which the points represent. For, in the first place, observed and recorded values may differ from true and actual values by reason of errors of measurement and of observation, and, secondly, there may be another variable quantity, unknown or ignored, like temperature in the example we have just considered, which contributes a secondary cause of the variable considered as an effect. If a number of plotted points lie approximately on a graph of simple mathematical shape, then the divergencies of the point positions from those of true correspondence with the graph are often assumed to

be due to what are called errors of observation. The drawing of the graph of a simple function that approximates most closely to actually containing the points is called fitting the graph, and this is an important technique which is described in the book on practical engineering mathematics and which is touched upon in Chapter XI. Graph fitting is really legitimate only when the kind of mathematical function connecting the variation of the two physical quantities is known on theoretical grounds. A higher technique in the interpretation of the graphical representation of observations is the detection and investigation of previously unsuspected secondary causes of variation.

It is not always, however, that the kind of mathematical function connecting two physical variables is known on the grounds of theory. Vapour pressure, for instance, is certainly a function of temperature, and magnetic flux is a function of magnetising force, but the underlying connections of these sets of physical variables is not known sufficiently well for the functional relationship to be expressed in a mathematical form. A graph showing simultaneous variation of this kind may be considered to be that of an unknown function. Every empirical graph is based upon a limited number of plotted points corresponding to a limited number of pairs of observations, and one of the practical uses of such a graph is to enable an unlimited number of further sets of simultaneous values to be inferred from the position of points on the graph lying between those corresponding to actual observations. Thus the graph of Fig. 43 is a smooth curve drawn approximately to contain all the points, indicated by small crosses, which correspond to observed values, and from this graph we can deduce another pair of simultaneous values of the variables from the co-ordinates of the point A. This is called interpolation, and its legitimacy depends evidently on the continuity of the function represented by the curve in the interval lying between the observed values on each side of A. When a pair of values of the variables is inferred from the co-ordinates of a pair of points on a portion of the graph produced beyond the range of the observed values, this is called extrapolation, and it is evidently much less reliable than inter-

polation, because the approximate form of part of a graph gives no certain indication of what lies beyond this range.

Sometimes an unknown functional relationship of two variables can, over a limited range of variation, be expressed very approximately by a simple mathematical function. Thus, at constant temperature, the pressure and volume of a gas are connected by the relation of inverse proportionality which is represented by a rectangular hyperbola, provided

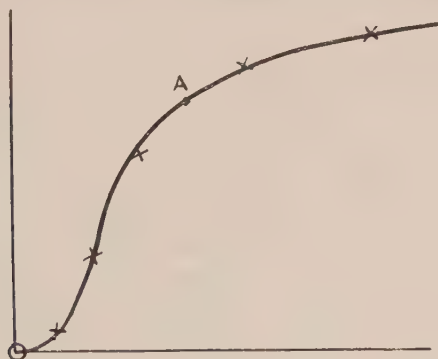


FIG. 43.

that the pressure is much less than its critical value. A physical phenomena can often be studied in an approximate way over a limited range by expressing an empirical function in a nearly equivalent mathematical form. Thus the unknown function represented by the graph of Fig. 43 can be approximately represented in this way. For, we see that for small values of x the graph seems like the parabola of Fig. 31, but that for large x values the graph might be asymptotic to a straight line parallel to the horizontal x axis. A function $y = \frac{ax^2}{b^2 + x^2}$ will therefore represent the graph approximately, for, when x is small and its square is negligible relatively to b^2 , the function is nearly the same as $y = \frac{ax^2}{b^2}$ which is a parabola, while, when x is so great that a^2 is negligible relatively

to b^2 the function approximates to $y = \frac{ax^2}{x^2}$ or $y = a$, a horizontal straight asymptote. The invention of mathematical functions approximately to correspond with empirical functions is another branch of technique of, perhaps, greater importance to the engineer than that of tracing the graphs of specified functions. As, however, this book is concerned with principles rather than technique we shall content ourselves with this short account of the matter.

CHAPTER VII

EQUATIONS

Basic Ideas. The statement that y is a function of x , symbolically expressed as $y=f(x)$, means that to every value of x there correspond one or more values of y . If the function is algebraic, or can be expressed in the form of an infinite converging series, then a numerical value of y corresponding to an assigned numerical value of x can, if the constants of the function are also given numerical values, be worked out by arithmetic. The converse problem of finding the value of x which will make $f(x)$ equal to an assigned value of y is much more difficult. The statement that $f(x)=a$, a constant, is called an equation, and the process of finding the value or values of x which make the statement true is called solving the equation. $f(x)=a$ is the same as $f(x)-a=0$, and as $f(x)-a$ is another function of x , an equation is usually stated in the general symbolic form $f(x)=0$.

If $f(x)$ contains, in addition to x symbols, only numbers, then the solutions give x in numerical form. Thus the equation $x-3=0$ has an obvious solution $x=3$, for 3 subtracted from 3 leaves zero. Again, the equation $x^2-5x+4=0$ can easily be seen to be true if x is equal either to $+4$, or $+1$, for taking the first value we find 4^2-20+4 , evaluated by arithmetic, is zero. If, however, the constants or coefficients in the function are literal (letters), as in $x^2-x(a+b)+ab=0$, the meaning of a solution is not quite so clear. If, however, we change x into $+a$ we get $a^2-a^2-ab+ab=0$, and this must be true whatever the values of a and b may be. A solution of an equation which can be obtained as a number is called a numerical solution, a solution of an equation expressed in more general terms with letter or literal coefficients gives x not as a number pure and simple, but as an expression or function containing one or more of the coefficients. This is called a formal solution. When a numerical solution is obtained and x in the equation is everywhere changed to this solution value, the arithmetical sum in the evaluation of

$f(x)$ gives 0 for an answer. If x is changed into a formal solution, an equation like $a^2 - a^2 - ab + ab$ is obtained which is true for all values of the symbols in it. Such an equation is called an identity. If $A=B$ is a symbolic statement of an identity, which is unconditionally true, then B must be merely another way of expressing A .

Equations are designated according to the kind of function of x that is equated to zero. When, as in the equation $a^x + x + b = 0$, a transcendental function a^x appears, the equation is called transcendental. If the function is algebraic the equation is algebraic, and an algebraic equation can always be made to depend upon the statement that an integral function of x , as defined on p. 92, is equal to zero. Such an equation is called an integral equation, or, often, simply an equation, without any qualifying term. The degree of the equation is the index of the highest power of x appearing in the function. If the highest power is 1, as in the equation $ax + b = 0$, the equation is called simple. The names quadratic, cubic, quartic, and quintic are applied to integral equations, the highest powers of x appearing in which are respectively, x^2 , x^3 , x^4 , and x^5 .

If an equation with numerical coefficients can be satisfied when x is a real number then a solution can always be obtained to as high a degree of approximation as may be required if the equation is integral and, usually, if it is algebraic or transcendental. If, however, the coefficients of the equation are literal, formal solutions for integral equations can be always obtained only if the highest power of x does not exceed 4. There is no general way of solving transcendental equations with literal coefficients. The difference between a formal and a numerical solution of an equation can be illustrated by the simple case $x^2 = a$. The formal solution is obvious, $x = \pm \sqrt{a}$, but if $a = 2$ a numerical solution can be obtained only approximately and only by a special arithmetical calculation. Again, no formal solution of the equation $x^5 + ax^2 + bx - c = 0$ can be obtained, but if a and b are each 1, and c is 38, a solution $x = 2$ is easily found by guessing or trial.

We have already considered some very easy examples of the solution of equations in our brief study of the technique

of the manipulation of engineering formulae on p. 87. Although this kind of technique is sufficient for most ordinary engineering calculations, the theory of equations is an important branch of mathematics the rudiments of which the engineer ought to understand. While, therefore, the subject-matter of this chapter may appear to bear less on specifically engineering mathematics than those which precede or follow it, and to be of a rather theoretical character, the reader will be well repaid for the mental effort required to grasp the rudimentary principles with which we now proceed to deal.

Simple Equations. Numerical simple equations are merely arithmetical sums expressed in algebraic language, and the solution of them is rarely difficult. Thus, considering $\frac{x-2}{3} + x = 22$, we first remove the fraction by multiplying the equal sides of the equation by 3, and, as 3 times $\frac{1}{3}$ of $(x-2)$ is simply $x-2$, we have the equivalent equation $x-2+3x=66$, or $4x=68$, which gives the solution $x=17$. That this is the solution can be checked by arithmetic, for $\frac{1}{3}(17-2)+17$ is evidently equal to 22. All this is very obvious, but the question should be asked, is 17 the only solution? or, is there any number other than 17 which substituted for x will make the equation true? The following paragraph will provide an answer to this question.

In however complicated form a simple equation may be stated, it is fairly evident that, as two classes of terms only are contained in the equation, x terms and constant terms, by collecting these classes together we can arrive at a modified equation in the standard form $ax+b=0$. Thus, for instance, the equation in the last paragraph was transformed or reduced finally to $4x-68=0$. The genesis of the equation $ax+b=0$ is the function $y=ax+b$, and this is made into an equation by making $y=0$. But, as we have seen on p. 111, this function is represented by a straight line graph with a slope a , and which cuts the vertical or y axis at a distance of b from the origin. The point on the graph that makes $y=0$ is that point where it intersects the x or horizontal axis, and the solution of the equation is the x co-ordinate of this point, or what we have

called the intercept of the graph on the x axis. As the graph can only cut the x axis once, there is only one value of x which will make $y = ax + b$ equal to zero, so that the equation $ax + b = 0$ can have only one solution.

Simultaneous Equations. The expression $ax + by + c$ has been defined on p. 85 as a function of two variables x and y , and if this is equated to 0, as $ax + by + c = 0$ we have y as an implicit function of x , which we can easily convert to the explicit form if we please, and which can be represented by a straight-line graph. Every point on the graph satisfies or fits the equation. Suppose we have a second equation $a_1x + b_1y + c_1 = 0$, similar in form, but with different coefficients. This will correspond to another straight-line graph, and the co-ordinates of every point on this graph will fit the second equation. Can there be found any pair of values of x and y that will fit or satisfy both equations? It is evident that this question is the same as, is there any point which is contained in the graphs of both equations? The answer is obvious; if the graphs intersect, the point of intersection is the one required, and the co-ordinates of this point will fit each equation as they refer to a point common to the corresponding graphs. But can such a point always be found? It can, provided the two graphs are not parallel.

Let us consider $ax + by + c$. We can convert this to an explicit function of y ; $y = -\frac{b}{a}x - \frac{c}{a}$, and this shows that the slope of the graph is $-\frac{b}{a}$. Similarly the slope of the graph representing $a_1x + b_1y + c_1 = 0$ is $-\frac{b_1}{a_1}$. If these slopes are equal the graphs are parallel, and no intersecting point can be found. Otherwise the graphs will meet, and the intersecting point will determine a set of values of x and y that will fit both equations, and as the graph intersect in one point only, this is the only fitting set that can be obtained.

A pair of equations like $ax + by + c = 0$ and $a_1x + b_1y + c_1 = 0$, taken together are said to be simultaneous, and, subject to the condition that the corresponding graphs are not parallel, a set of values of x and y can always be found that will fit

both equations. One way of finding these values is fairly obvious, for if we convert both equations into explicit functions of y then the solution for x is contained implicitly in the statement that the y of one function is equal to the y of the other. Thus if one equation is such that it can be converted to $y=2-x$, and the other to $y=8-2x$, we have, equating the y 's, $2-x=8-2x$, which is the same as $x=6$. This value of x ought evidently to give the same value of y , in the two x functions representing the equations; $y=2-x$, gives for $x=6$, $y=-4$, and it is seen at once that $y=8-2x$ gives the same y value for $x=6$. $x=6$ and $y=-4$ are the solutions of the two simultaneous equations, and the simple equation containing x only. $2-x=8-2x$ is obtained by what is called eliminating y .

Let us consider $y=2x$ and $y=4x$ as a pair of simultaneous equations. Proceeding as before we obtain $2x=4x$. How can twice x equal four times x ? The answer is that it can only do so if $x=0$. This gives $y=0$, for $y=2 \times 0$ from one equation. That this is the solution is evident, for both the graphs corresponding to the two functions of y pass through the origin point and intersect there. Again, consider $y=2x+8$ and $y=2x-9$. Equating y 's gives us $2x+8=2x-9$, or subtracting $2x$ from each side of this equation $8=-9$, an impossible result. This means, of course, that there is no solution to the equations in the ordinary sense of the word, and the reader ought to see at once that the graphs representing the equations are parallel, each having a slope of 1. Mathematicians interested in theory do not, however, like these anomalous cases, and they say that the solution to a set of equations like this is that both x and y are infinite, an algebraic equivalent to the statement that the intersecting point of the parallel graphs is at an infinite distance from the origin point. $2x+8=2x-9$ means that x must be so large that the results of adding 4 to it and subtracting $4\frac{1}{2}$ from it are indistinguishable. If this idea of infinite solutions seems far-fetched and artificial to the reader, he need not trouble himself about it. For all practical purposes, we can say that no solution can be obtained of the equations $y=2x+8$ and $y=2x-9$.

Let us now return to the pair of equations $ax + by + c = 0$ and $a_1x + b_1y + c_1 = 0$. We can obtain the formal solution in the manner just demonstrated, that is, by turning each equation into an explicit function of x . If we proceed in this way we find the solution is

$$x = \frac{bc_1 - b_1c}{ab_1 - a_1b} \quad \text{and} \quad y = \frac{ca_1 - c_1a}{ab_1 - a_1b}$$

as the reader can work out for himself if he has the necessary interest and patience. These expressions are worth some study; the order of the letters of each product is always cyclic, that is, in the sequence $abcab$ and so on, the denominators of the two fractions are the same, and in both numerators the missing letter is the one associated with x or y as the case may be. If we denote the x numerator by D_1 , the y numerator by D_2 , and the common denominator by D_3 the solution is, symbolically, $x = \frac{D_1}{D_3}$, $y = \frac{D_2}{D_3}$ and the D_1 , D_2 and D_3 are often represented by the notation

$$D_1 = \begin{vmatrix} b & c \\ b_1 & c_1 \end{vmatrix} \quad D_2 = \begin{vmatrix} c & a \\ c_1 & a_1 \end{vmatrix} \quad D_3 = \begin{vmatrix} a & b \\ a_1 & b_1 \end{vmatrix}$$

These square collection of letter symbols are called determinants, and, in this case as they represent the difference of the products of two letter symbols, they are of the second order. The products represented by the determinant are those of diagonally opposite letters, the product going down to the left being subtracted from the other.

Now let us consider a third equation $a_2x + b_2y + c_2 = 0$ in connection with the previous two for which a solution can always be found provided $\frac{b}{a}$ is not equal to $\frac{b_1}{a_1}$. Can a set of values of x and y determined by the third equation be found to fit the other two? The geometrical equivalent of this question is, can a point be found on the line graph represented by $a_2x + b_2y + c_2 = 0$ which form is common to the line graphs represented by the graphs of the other two equations? The obvious answer is, only if the third graph passes through the intersection of the other two. If the third graph does not satisfy this condition no point on it will fit the other two

equations and the set of three are said to be inconsistent, that is, they cannot all be simultaneously true. If, however, the three straight line graphs all pass through one point then the co-ordinates of this point have x and y values which are a solution to any of the three pairs of equations. The set is

then said to be consistent. To put this symbolically, $x = \frac{D_1}{D_3}$, $y = \frac{D_2}{D_3}$ are the solution of $ax + by + c = 0$ and $a_1x + b_1y + c_1 = 0$.

If $a_2x + b_2y + c_2$ is consistent with the other two equations, then this solution will fit this last equation. Otherwise

$\frac{D_1}{D_3} \times a_2 + \frac{D_2}{D_3} \times b_2 + c_2 = 0$ will be true. Multiplying every term in this equation by D_3 we find, as $D_3 \times 0 = 0$, that $a_2D_1 + b_2D_2 + c_2D_3 = 0$, and that this is the condition that the three equations are consistent. This statement is generally exhibited in the form

$$\begin{vmatrix} a & b & c \\ a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \end{vmatrix} = 0.$$

The square arrangement of letters forming a determinant of the third order. This, as appears from the previous equation, is obtained by associating each of the symbols a_2 , b_2 , and c_2 with the second-order determinants formed by the letters in the rows above not contained in the column containing the multiplying letter. Thus a_2 in the first column, is multiplied by the determinant formed by letters in the upper two rows and the second and third columns.

The theory and the technique of determinants is a distinct branch of mathematics of considerable importance in its higher branches, and in some special kinds of engineering calculations. We have merely introduced the determinant notation in order that the reader may know in a general way what it implies. One of the most useful ways of regarding a determinant of the third order as given above is as a concise statement of the relation that must exist between the co-efficients of a set of three equations containing two unknowns or variables, x and y , in order that the equations may be consistent.

We can have a set of three simple equations containing three unknowns, x , y , and z , and to these one unique set of values or a solution can generally be found that will fit all three equations. As a general rule it may be said that provided the number of unknown quantities is the same as the number of equations all considered simultaneously true, a solution giving the value of the unknowns can usually be found. The subject of the solution of sets of three or more simultaneous equations is a matter rather of the technique than of the basic theory of mathematics, and for further information on this subject the reader may, if he requires it, consult a work on applied mathematics. He will find there that the solution in the symbolic determinant form can be written down by a rule precisely similar to that applicable to a set of two equations, but he will soon see that the evaluation of the determinants of a higher order than the third is a technique calling for considerable skill and experience.

General Remarks on Integral Equations. When a simple equation containing no power of x higher than the first is converted to the general form $ax + b = 0$ it gives its solution almost explicitly. This, however, is not the case for equations of higher degrees, because the integral function equated to zero generally contains all powers of x up to the highest, the index of which determines the degree of the equation. Thus a quadratic equation can always be converted to the standard form $ax^2 + bx + c = 0$ and there is here no obvious way of obtaining the value or values of x which satisfy the equation. The problem of the solution of equations of a higher degree than the first depends upon a somewhat difficult branch of algebraic theory, a brief and simple sketch of which we shall now give.

In Chapter II we have encountered some very simple examples of factors of an integral function; we found, for instance, on p. 14, that the integral function $x^2 - a^2$ is equal to the product of the simpler functions $(x - a)$ and $(x + a)$, and this means that if $x - a$ is multiplied first by x and then by a , the two terms xa destroy each other, and the result is $(x^2 - a^2)$. $x^2 - a^2 = (x - a)(x + a)$ is equivalent to saying that $(x - a)$ and $(x + a)$ are the factors of $x^2 - a^2$. We must be very careful to

understand fully what this means. Consider the function $x^3 + a^2$; now it is impossible to find any two integral functions with real number terms that, multiplied together, algebraically, will give a result $(x^2 + a^2)$. In fact $(x^2 + a^2)$ is something like a prime number, but not altogether. If, for instance, $x=3$ and $a=1$ the function is equal to $3^2 + 1 = 10$, and it has the number factors 5 and 2; if, however, $x=4$ and $a=1$, then the function is equal to 17, a prime number. Numerically, $x^2 + a^2$ gives a composite or factorisable number only for certain values of x and a . The function $x^2 - a^2$ is different; assigning whole number values to x and a , the function is always equal to a composite number, having the factors $(x-a)$ and $(x+a)$.

Let us consider the equation $f(x)=0$ where $f(x)$ is integral, and consists of the sum of powers of x each multiplied by a constant or coefficient, together with a coefficient not associated with x , and often called an absolute term. If $(x-a)$ is a factor of $f(x)$, then $f(x)=(x-a) \times \psi(x)$ where $\psi(x)$ is another integral function of a degree one lower than $f(x)$, because in working out the product the term of the highest power of x in $\psi(x)$ will be multiplied by x and have the index of its power raised by 1. If $(x-a)$ is not a factor of $f(x)$ it will always be possible to find a function $\psi(x)$ one degree lower than $f(x)$ such that $f(x)=(x-a) \times \psi(x) + R$ where R is a quantity independent altogether of x , that is, containing no x terms. R is like the remainder in a division sum, and it is called a remainder in algebra. The statement given is true for all values of x , and so it is true when $x=a$. Then $f(a)=(a-a)\psi(a) + R$ by making x equal to a . But as $a-a=0$, the product $(a-a)\psi(a)$ is also 0, so that $R=f(a)$. If therefore we follow up the arithmetical analogy we can say that dividing $f(x)$ by $(x-a)$ gives a remainder which is the same function of a as $f(x)$ is of x . This is known in algebra as the Remainder Theorem. If $R=0$, then $f(a)=(a-a) \times \psi(a)=0$, and this satisfies the condition of the integral equation $f(x)=0$. a is a quantity that satisfies the equation, or, as it is sometimes expressed, which causes $f(x)$ to vanish. It is evidently a solution or root of the equation. We thus arrive at the important conclusion that if $(x-a)$ is a factor of $f(x)$ then one

root of the equation $f(x)=0$ is found by putting $x-a=0$, and its converse, that if a causes $f(x)$ to vanish then $(x-a)$ is a factor of $f(x)$. To illustrate this we may consider the equation $x^3-x^2+x-1=0$. We see at once that the function vanishes when $x=1$, for it is then $1-1+1-1=0$. $(x-1)$ is therefore a factor of x^3-x^2+x-1 , as can be seen at once by converting the function to $x^2(x-1)+1 \times (x-1)$ or $(x^2+1)(x-1)$.

Let us take another kind of illustration. As the equation $x^2-a^2=0$ is the same as $x^2=a^2$, the solutions or roots are seen at once to be $\pm a$. But the original equation can be written $(x-a)(x+a)=0$. As $(x-a)$ is a factor of the function equated to zero, one root is obtained by writing $x-a=0$, this is $x=a$. But $(x+a)$ is also a factor, so another root is obtained from $x+a=0$, and this is $x=-a$.

It is a fundamental theorem of algebra that every integral equation has at least one root, that is, that there is at least one value of the unknown quantity or variable which causes the function to vanish. We must take this for granted, not because it is self-evident, but because a formal proof is difficult and beyond the scope of this book. We note also that the statement is that there is at least one root, not that this root can necessarily be found in any particular way. Taking this statement as true it means that, as $f(x)=0$ has one root at least, which may be real, imaginary or complex, then signifying this root by a , $(x-a)$ is a factor of $f(x)$ so that $f(x)=(x-a) \times \psi(x)$. But $\psi(x)$ is an integral function of x of one degree lower than $f(x)$, $\psi(x)$ must therefore have at least one root, say b , so that it can be put $\psi(x)=(x-b) \times \phi(x)$ where $\phi(x)$ is an integral function two degrees lower than $f(x)$. Proceeding in this way we see that $f(x)=0$ must have a number of roots equal to the degree of the function, that is, equal to the index of the highest power of x in it, so that if the roots of the equation are a, b, c and so on $f(x)$ can be expressed in the form $f(x)=p(x-a)(x-b)(x-c) \dots$, the number of factors being equal to the degree of $f(x)$ and p being a quantity independent of x . It will be evident, with a little thought, that as the highest power of x in $f(x)$ is obtained by multiplying together p and all the x 's in the factors, p is the coefficient of this highest power.

But can there be more roots than correspond to the index of this highest power? Suppose there is an additional root, s , then, substituting this for x in the function we have $p(s-a)(s-b)(s-c) \dots = 0$; and as s is different from all the other roots, none of the factors are zero, so we are forced to the conclusion that the quantity p must be zero, and the function cannot be of the degree we have assigned to it. But proceeding in this way with what is left after rejecting the term having the highest power of x , we can show similarly that the coefficients of all the terms must be zero. Thus, if the equation $f(x)=0$ is satisfied by more values of x than the index of the highest power in the function, the coefficient of every power and the absolute term are all zero, so that the equation is satisfied by every value of x . Such an equation is called an identity as has already been explained briefly on p. 134, but which we shall illustrate in a little more detail. Consider $x^2 - 9 - x(x+3) + 3(x+3) = 0$. This looks like a *bona fide* equation, the highest power of x being 2, and which ought to have two roots. If we put $x=1$ the equation becomes $1 - 9 - 4 + 12 = 0$, which is right; also, putting $x=2$, we have $4 - 9 - 10 + 15 = 0$, which is also right, so that $x=1$ and $x=2$ might be considered to be the roots. But if we put $x=3$ we get $9 - 9 - 18 + 18 = 0$, which is also right, so that the equation is satisfied by more than two values of x . It must therefore be an identity. We soon see that this is so by carrying out the multiplications of the bracketed quantities for $x^2 - 9 - x(x+3) + 3(x+3)$ is $x^2 - 9 - x^2 - 3x + 3x + 9$ or $x^2 - x^2 + 3x - 3x + 9 - 9$, which is zero, so that the equation is simply $0=0$, and is true for every numerical value we care to assign to x in its original form.

The proposition that if the integral equation $f(x)=0$ is an identity, then the coefficients of all powers of x are zero leads to a very important and far-reaching corollary. Suppose that a function of x is equal for all values of the variable to the series of a number of powers of x , or, that $f(x) = ax^n + bx^{n-1} + cx^{n-2} \dots$ for all values of x , that is, suppose that the power series is a legitimate way of representing the function algebraically. Now suppose that by some process of mathematical reasoning we arrive at the result that the same function

of x is equal to an apparently different power series $Ax^n + Bx^{n-1} + Cx^{n-2} \dots$ also for all values of x . It then follows that the two power series are always equal and

$$ax^n + bx^{n-1} + cx^{n-2} \dots = Ax^n + Bx^{n-1} + Cx^{n-2}.$$

So that

$$(A - a)x^n + (B - b)x^{n-1} + (C - c)x^{n-2} \dots = 0$$

is an equation true for every value of x . This equation is an identity, and the coefficients of every power of x , and also the terms, if any, not containing x , must be all equal to zero. Thus $(A - a) = 0$, and $A = a$; similarly $B = b$, $C = c$, and so on. In other words, in the apparently different series of x powers representing the same function the coefficients of like powers of x are equal.

This principle, commonly called that of undetermined coefficients, is so important in mathematics that we shall illustrate it by a concrete example. Suppose that a function of x can be expressed in the form of a power series, as

$$f(x) = 1 + ax + bx^2 + cx^3 \dots$$

this meaning that for suitable values of the coefficients a , b , c , and so on, the series is always equal to the function. Now suppose that by some process of reasoning we find that the same function can be expressed as

$$f(x) = a + 2bx + 3cx^2 + 4dx^3 \dots$$

where the a in this second expression is the same as the a in the first. Then the principle of undetermined coefficients is that in the two power series the coefficients of similar powers of x are equal. By this principle we find that $a = 1$, for a and 1 may be considered coefficients of the zero-th power of x . $2b = a = 1$ so that $b = \frac{1}{2}$. $3c = b = \frac{1}{2}$, so that $c = \frac{1}{2 \times 3}$ and

so on, the coefficient of the term next to the n th being the n th term coefficient divided by n . The two power series with undetermined coefficients has therefore enabled us to arrive at the function already mentioned on p. 98.

$$1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} \dots$$

in which the coefficients are actual and determined fractional numbers. These coefficients it must be noted have been evaluated because by a process of reasoning we have been able to state the function with undetermined coefficients in two forms. We shall see later how a second statement, such as we have assumed possible, can be obtained. The important thing at the moment is for the reader to understand the reasoning that has led up to the principle of undetermined coefficients, so that he will be able to give it conscientious mental assent. The principle is often assumed in books on practical mathematics to be almost self-evidently true, which it certainly is not.

Let us now return to the statement that an integral function of x of the n th degree can be put in the form

$$f(x) = p(x-a)(x-b)(x-c) \dots$$

where the number of factors is equal to the index of the highest power of x in the function, and a, b, c and so on are the roots of the equation $f(x) = 0$. If all the coefficients of x powers in the function are real numbers it does not follow that all the roots will be real. The function $x^2 + a^2$ where a is real has no real factors, but as $x^2 + a^2 = 0$ can be turned into $x^2 = -1 \times a^2$ it is seen at once that $x = \pm ja$ and the equation has two imaginary roots. It is not difficult to see that when an equation has real coefficients, then imaginary or complex roots must occur in conjugate pairs so that the product of the factors containing them may be real. Again, suppose that a function, like $x^2 - 2ab + b^2$, can be expressed as the product of equal factors $(x-a)(x-a)$. The equation $(x-a)(x-a) = 0$ in one manner of speaking has one root only, a , although being of the second degree it ought by general theory to have two. We shall refer again to this point in the following section, but we may say here that, as with the anomalous case mentioned on p. 137, mathematicians do not like to acknowledge exceptions to a general rule, so that they say that the equation $(x-a)(x-a)$ has two equal roots.

Graphical Illustrations. If the graph of the function $f(x)$ is plotted, then the roots of the equation $f(x) = 0$ or the values of x which make the function take a zero value are given by

the x co-ordinates of the points where the graph intersects the horizontal axis where $y=0$. Thus the function whose graph is as shown in Fig. 44 determines two roots x_1 and x_2 of the equation obtained by putting this function equal to zero. If the graph of a function were like Fig. 44 but its maximum point came just above the

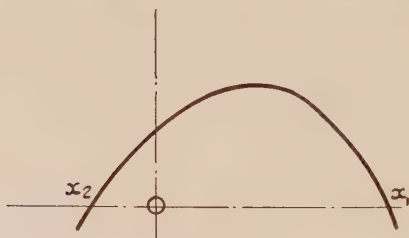


FIG. 44.

horizontal axis as in Fig. 45 the two roots x_1 and x_2 of the corresponding equation would be very nearly equal, and when the graph was in such a position that the x axis was tangential to it the two co-ordinates determining the roots of the equation would merge into one. When therefore an equation has two equal roots this means that the graph of

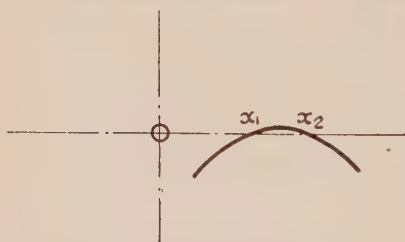


FIG. 45.

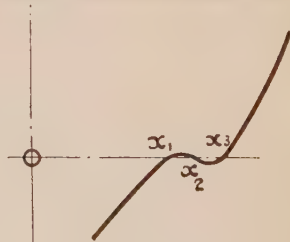


FIG. 46.

the function has the x axis for a tangent at the point whose abscissa gives these roots.

It is possible to have more than two equal roots. The graph in Fig. 46 cuts the x axis in three points quite near to each other, and if with such a graph the three intersecting points merge into one we shall have three equal roots of the corresponding equation. In this case the x axis will be a tangent intersecting the graph.

The equation $(x-a)(x-a)(x-a) \dots = 0$ has as many equal roots as it has equal factors. The equal roots determine a point where the x axis is tangential to the graph of the function and the number of the roots determine what is called the order of contact. If this order is odd the tangential axis of co-ordinates crosses the graph.

The graph of an integral function will not necessarily cross or touch the x axis. Fig. 47, full line graph, is that of the function $y=x^2+1$, and this graph lies entirely above the x axis. We know from algebraic consideration that the roots of the corresponding equation $x^2+1=0$ are $\pm\sqrt{-1}$ or $\pm j$; but what geometrical meaning can be found for this state-

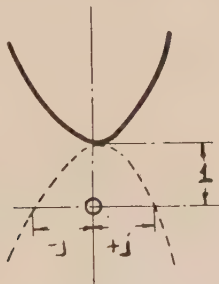


FIG. 47.

ment? If imaginary values are assigned to x in the function, the x ordinates corresponding in the geometrical sense will be measured in a direction at right-angles to the plane of the paper and their squares will be negative, so that the function in a perpendicular plane becomes $y = 1 - x^2$. The graph of this latter function is shown dotted; it is the graph for imaginary x values turned about the y axis through 90 degrees into the plane of the paper. The graph intersects the imaginary axis of x at points $+1$ and -1 so that the roots of the original equation are $+j$ and $-j$.

Suppose we can find two values x_1 and x_2 of variable in $f(x)$ which, as shown in Fig. 48, give small values y_1 and y_2 of the function, the one positive and the other negative. If the function is continuous in the interval between these values of y then the graph must cross the x axis in this interval and

$f(x)=0$ must have a root between x_1 and x_2 . If the interval is small, so that the portion of the graph corresponding is approximately straight, the graph will determine two triangles, shown shaded in the figure, which are of similar shape; and if x is the root or the co-ordinate where the graph crosses the axis then $\frac{x_1 - x}{x - x_2} = \frac{y_1}{y_2}$ by the principle of the proportionality of corresponding sides of similarly shaped triangles. As in this equation x_1 , x_2 , y_1 and y_2 are known, it is a formula whereby an approximate value of the root x may be estimated.

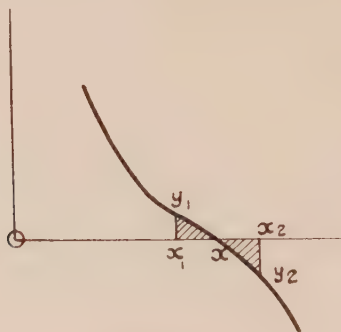


FIG. 48.

Again, suppose that assigning a value x_1 to $f(x)=y$, we find a very small value corresponding, y_1 . Then if we know the slope of the graph at x we have all the particulars of the left-hand shaded triangle in Fig. 48, and as the slope is $-\frac{y_1}{x_1 - x}$ if y_1 is positive, a formula can be obtained for finding x , the root, on the assumption that in the interval x_1 to x the graph can be considered to be nearly a straight line. This is the basis of Newton's method for the approximate solution of equations, and it is applicable to any equation, provided a value x_1 can be found giving a small value of the function, that the slope at x can be calculated, and that the function is continuous from x_1 to the intersecting position determining the root.

Quadratic Equations. The general form of an equation of

the second degree or of a quadratic equation is $ax^2 + bx + c = 0$ where a , b , and c are constants. By dividing every term in this equation by a it can, as $0 \div a = 0$, be converted to the form $x^2 + Bx + C = 0$ where $B = \frac{b}{a}$ and $C = \frac{c}{a}$. A quadratic equation has two roots so that, calling these α and β , $x^2 + Bx + C = (x - \alpha)(x - \beta) = x^2 - x(\alpha + \beta) + \alpha\beta$ for every value of x . By the principle of undetermined coefficients, therefore, $B = -(\alpha + \beta)$ and $C = \alpha\beta$, or in words, the x coefficient is minus the sum of the roots and the absolute term which does not contain x is equal to their product.

This principle enables us to find the roots of some quadratic equations by inspection. Thus, consider $x^2 - 5x + 4 = 0$, $-(-5)$ or $+5$ is the sum of the roots and 4 is their product. It is seen at once that the roots are $+4$ and $+1$, so that the equation in the same as $(x - 4)(x - 1) = 0$. Let us take a more difficult example: $x^2 + 9x - 112 = 0$. Here the sum of the roots is -9 , and their product is -112 . One root must therefore be positive and the other negative; as the sum is negative the negative root is numerically the greater, so that 9 is a numerical difference. A little thought and mental calculation soon shows that the roots are -16 and $+7$, so that the equation can be written $(x + 16)(x - 7) = 0$. This way of solving quadratic equations is sometimes described as by factorising, but the roots have really to be found before the factors can be written down. Factorising methods figure largely in school textbooks on algebra; they require considerable ingenuity and give good practice in mental arithmetic, but have little to do with real mathematical technique, first, because these so-called methods of solution are really processes of obtaining roots by inspection and mental calculation, and, secondly, because only a limited number of special quadratic functions can be factorised.

Let us, for example, take the equation $x^2 + 2x - 1 = 0$. We cannot find a pair of whole numbers that give a sum of -2 and a product of -1 , and a less direct method of solving the equation is evidently necessary. To develop a method we observe that $x^2 + 2x$ is part of the square quadratic function $x^2 + 2x + 1 = (x + 1)^2$ so that if we add to and subtract from the

left-hand side of the equation the same quantity 1 and so leave its value unchanged we shall get $x^2 + 2x + 1 - 1 - 1 = 0$, or $(x+1)^2 - 2 = 0$, or $(x+1)^2 = 2$ from which we obtain $x+1 = \pm\sqrt{2}$; is a combination of two simple equations from which we deduce two roots $x = -1 - \sqrt{2}$, and $x = -1 + \sqrt{2}$. Again, take the example $x^2 - 4x + 8 = 0$. $x^2 - 4x$ is part of the square $x^2 - 4x + 4 = (x-2)^2$, so that, adding and subtracting 4 we get $x^2 - 4x + 4 - 4 + 8 = (x-2)^2 + 4 = 0$, whence $(x-2)^2 = \pm\sqrt{(-4)}$ which square root can be written $\pm 2j$. Hence we obtain $x-2 = \pm 2j$, and the two solutions $x = 2 + 2j$ and $x = 2 - 2j$ are obtained, a pair of conjugate complex numbers.

This artifice can be applied to any kind of quadratic equation. Take the most general form $x^2 + \frac{b}{a}x + \frac{c}{a}$, how can we turn $x^2 + \frac{b}{a}x$ into the form of the square of a binomial like $(x+p)^2$; a little thought will show that as $\frac{b}{a}$ is twice the other term of the binomial, the term to be added and subtracted is $\left(\frac{b}{2a}\right)^2$ so that the equation takes the form $x^2 + \frac{b}{a}x + \left(\frac{b}{2a}\right)^2 - \left(\frac{b}{2a}\right)^2 + \frac{c}{a}$, or $\left(x + \frac{b}{2a}\right)^2 = \left(\frac{b}{2a}\right)^2 - \frac{c}{a}$, and by taking the square root of each side we can obtain two solutions in which x is a function of a , b , and c . We see at once that if $\frac{c}{a}$ is greater than $\left(\frac{b}{2a}\right)^2$ the square of $\left(x + \frac{b}{2a}\right)$ is a negative number and the roots are complex numbers, which are conjugate. If, $\left(\frac{b}{2a}\right)^2 - \frac{c}{a}$, when evaluated by assigning numerical values to a , b , and c , is a square number, the roots will be rational, but not otherwise. If $\left(\frac{b}{2a}\right)^2 - \frac{c}{a}$ is equal to zero the equation becomes $\left(x + \frac{b}{2a}\right)^2 = 0$, which is $\left(x + \frac{b}{2a}\right)\left(x + \frac{b}{2a}\right) = 0$, and which shows that the two roots $-\frac{b}{2a}$ are equal. A manipulation of

the last form of the quadratic equation leads to a general formula for the solution of $ax^2 + bx + c = 0$. This formula is

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}.$$

This formula is easily forgotten, and, as it can always be found in any mathematical textbook if required, there is no need to try to remember. What is important for those who wish to understand something of what mathematics mean, is that the algebraic process of obtaining the formula should be grasped and understood. This process is essentially that of converting a quadratic function $ax^2 + bx + c$ into the difference of two squares, and if the reader pursues his studies he will find that this technique is of considerable importance in branches of mathematics higher than that of the solution of quadratic equations.

Formal Solution of Higher Equations. The reader who has grasped the principle whereby a formal solution of a quadratic equation has been obtained will see that this depends upon converting the equation $ax^2 + bx + c = 0$ to another equivalent equation $(y^2 - A^2) = 0$, where the new variable or unknown y depends, according to a definite rule, on the original one x .

By a somewhat analogous but much more complicated process an equation of the third degree, or a cubic, expressed generally as $ax^3 + bx^2 + cx + d = 0$ can be converted into the form $y^6 + Ay^3 + B = 0$, or, what is the same, $(y^3)^2 + A(y^3) + B = 0$ where y depends upon a definite way upon x . This latter equation can be solved as a quadratic, of which the unknown is (y^3) , and hence the solutions of the original equation obtained. Since the quadratic has two roots, and each (y^3) solution has three cube roots, this method appears to give six roots of the original cubic, but it can be shown by algebra that these six roots really consist of three pairs, each pair being identical. The best-known method of solution of this kind is known as Cardan's method.

The formal algebraic solution of a cubic leads to an interesting and anomalous result, real solutions come out as the sum of the cube roots of two conjugate complex numbers, as an

expression like $\sqrt[3]{a+jb} + \sqrt[3]{a-jb}$. This puzzled mathematicians before complex numbers were thoroughly understood, but the reader who has grasped the latter part of Chapter IV will see easily that $a+jb$ and $a-jb$ are represented geometrically by equal directed steps inclined at equal angles to the reference line, the one above and the other below, the cube roots of $a+jb$ and $a-jb$ are also equal steps inclined at one-third the equal angles, the one above and the other below, the reference line, so that the sum of the cube roots is represented by the resultant of these latter steps, and hence by a step in the reference line of real numbers.

By an even more complicated algebraic process a quartic equation, or one of the fourth degree can be made to depend upon a cubic, and hence can be solved formally. Much labour was at one time expended by mathematicians in trying to solve a quintic or fifth degree equation formally, by making it depend upon a quartic. It is now known that this cannot be done and that, as has been previously stated, no algebraic or formal solution can be obtained of an integral equation with literal coefficients if its degree is higher than four. Formal solutions of third and fourth degree equations are of little use in practical mathematics for the evaluation of numerical roots, and for this purpose other methods are generally used.

Approximate Solution of Equations. We have seen that if the whole graph of the function $f(x)$ can be drawn, the approximate roots of the equation $f(x)=0$ are given by the x co-ordinates of the points of intersection of the graph with the x axis. The evaluation of roots by this process is a method of interpolation, that is, assuming the shape of the graph between two points the co-ordinates of which have been actually calculated. The graphical solution of an equation can often be made easier by an artifice whereby, instead of the actual graph of the corresponding function, two simpler graphs are plotted. This artifice can be illustrated by the following example. The equation $x^3 - ax^2 - b = 0$ can be changed to $x - a - \frac{b}{x^2} = 0$, by dividing every term by x^2 , and

this is the same as $y - y_1 = 0$ where $y = -a + x$ and $y_1 = \frac{b}{x^2}$.

A root of the original equation, therefore, corresponds to the condition $y = y_1$, so that if the y and y_1 graphs are plotted the root corresponds to the x co-ordinate of the point where the y co-ordinates of the graphs are equal, that is, to the point where they intersect. This device is useful, because the graph $y_1 = \frac{b}{x^2}$ can be sketched rapidly with fair accuracy; the graph $y = -a + x$ is of course a 45 degree straight line cutting the vertical axis at a distance a below the origin point. The application of this method is illustrated in Fig. 49, where

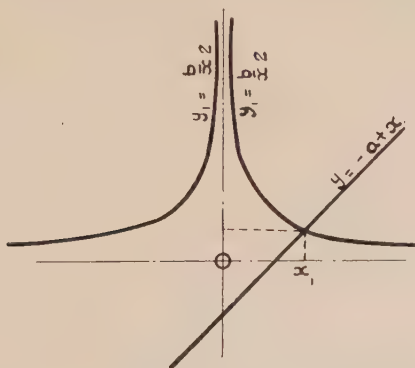


FIG. 49.

x_1 is the approximate root sought. This is one of the many devices for the approximate solution of equations for which experience and judgment are required in the application.

We have already referred on p. 148 to two methods whereby, knowing a root roughly, a more accurate value can be obtained by calculation, although in Newton's method we have not yet explained how the slope of the graph in the neighbourhood of the root value can be obtained. This point will be dealt with in Chapter VIII.

There is a standard method of obtaining numerical solutions to integral equations of any degree, which was invented by W. G. Horner at the early part of last century. The technique of this method is so powerful, and for those who can use it, enables high approximations to be so rapidly

obtained that, as Horner was largely a self-taught man, it has been rather ungenerously suggested that his remarkable invention was due to luck rather than to inborn mathematical talent of a high order. Horner's method is not difficult to learn, but, like most rather complicated mathematical techniques, it is very easily forgotten unless constantly used. As we are concerned in this book with basic mathematical ideas rather than with the details of calculating processes we shall merely, in concluding this chapter, indicate in a general way the underlying principle of Horner's method.

Suppose that the graph A in Fig. 50 is part of that of the integral function $f(x)$, and suppose that by actual calculation

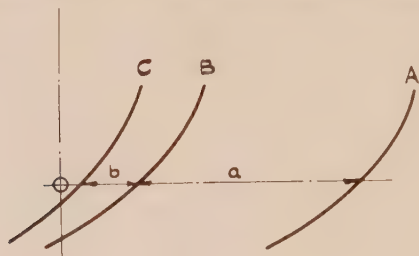


FIG. 50.

we find that the root of the equation is a little greater than some number a , that is, substituting a in $f(x)$ gives y a small positive value and substituting a number somewhat less than a gives y a small negative value. If we can construct another function of x , $\psi(x)$, having a graph exactly like A but moved a distance a to the left to the position B, this function will be of the same degree and kind as $f(x)$, but the value of x that makes it vanish will be a less than the root of $f(x)=0$. When we now try to get an approximate root of the new equation $\psi(x)=0$ where $\psi(x)$ corresponds to the graph B, the arithmetical work will be easier because we shall know that this root will be much less than those we used in the first trial. By this second trial we might find that the root of the equation $\psi(x)=0$ is a little greater than a number b . Making a third function corresponding to the graph moved to the position C, we may arrive at an equation

with so small a root that the position of this graph which intersects the axes of co-ordinates can be considered a straight line, so that all terms containing powers of x higher than the first in the function can be neglected and that its approximate root C can be found by solving a simple equation. The root of the original equation will then be the sum of the quantities a , b , and c . Horner's method is a technique whereby, first, the new functions corresponding to the leftward shift of a graph can be rapidly determined, and, secondly, whereby the best amounts of these leftward shifts can be easily estimated.

CHAPTER VIII

PERIODIC FUNCTIONS

Naturally occurring phenomena are seen by the most superficial observer to fall into two well-defined classes. The first, like day and night, the cycle of the seasons, and the appearance of the moon, in which the phenomena are regularly and continuously repeated, and the second, like the life and growth of living organisms and the changes in the surface of the earth by the destructive or constructive forces of nature, in which the efflux of time produces changes that are continuously progressive. Phenomena of the first class are usually called periodic, and those of the second class secular. These two words, periodic and secular, having a quantitative connotation are inherently mathematical in their basic significance, and the functional relationship between the measure of a phenomenon and the time interval reckoned from any instant is quite different in the two classes. The functions we have so far considered are inherently representative of secular changes. Thus, to consider the simplest of algebraic functions, that of direct proportionality, we are sure that if an effect is proportional in its magnitude to time efflux then, however small this effect may be in any assigned time unit, day, year, or century, the magnitude of this effect will increase continuously. In the case of inverse proportionality we have the reverse effect of a progressive and continuous reduction. Functions which represent the other class of natural phenomena are quite different in their nature, because efflux of time, or the continuous increase of the independent variable will continually repeat the conditions existing at any assigned instant. This leads to a simple definition of a periodic function in which time is the independent variable: $f(t) = f(t + mT)$, which is a symbolic statement of the fact that a quantity, say y , has exactly the same value at any time denoted by t that it has when t is increased by any exact multiple of a definite or fixed time interval T . T is the time interval, the elapse of which will always reproduce the value

of y , the function of variable time, t . T is here called the periodic time, and the number of T intervals occurring in any time unit, in a second, a minute, an hour and so on, is called the frequency per second, per minute, or per hour. Physical functions can be periodic in respect to distance; if the vertical undulations of a straight road are periodic, then the height of the road above a datum will be reproduced for equal movements of the distance measured along it. This could be expressed symbolically as $y=f(x)=f(x+m\lambda)$ where m is any whole number and λ is a constant increment of x . If a function, periodic in relation to any independent variable, is plotted into a graph the result will be like Fig. 51, in which

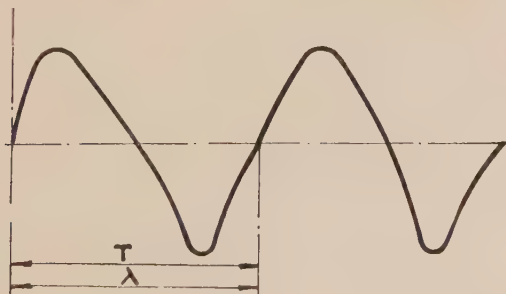


FIG. 51.

the curve of the graph continually repeats. The graph of a periodic function exhibits very clearly three basic characteristics of the function: first its period, in respect to time or length which is, graphically, the distance which includes a complete sample of the curve that is continuously repeated; secondly, the maximum variation of the dependent variable from the datum or zero value, which is known as the amplitude; and, thirdly, the manner in which the dependent variable changes during any period, which is shown by the shape of the graph, and which is commonly called the wave form.

It is possible to have a function simultaneously periodic in respect to two independent variables, space, and time. This is called a wave motion, and is illustrated by ripples or waves on the surface of a large area of water. At any instant of

time an instantaneous photograph of the surface of the water would be a graph like Fig. 51, periodic in relation to distance. At any point on the water surface the height of the surface is a periodic function of time. The first function has a period represented by a length λ , and the second by a time T , and the simultaneous variation in respect of space and time makes the maximum displacement of the water surface move by a distance λ in a time T , that is, with a speed equal to $\frac{\lambda}{T}$, which is called the velocity of propagation of the wave. The mathematical study of wave motion is difficult and beyond the scope of this book.

Periodic functions are of great importance in applied science. The regulation of a clock is governed by a pendulum, the movement of which is periodic. In modern electricity supply both current and pressure are continuously varying periodically with respect to time with a frequency, in this country, of 50 periods per second. There is a particular kind of periodic function which is not only basically the simplest in mathematics, but which is of the utmost importance because all other periodic functions have been shown to depend upon it. Objectively this function is represented very nearly by the movement of a very long pendulum or by the vibrations of the prong of a softly-sounding tuning fork. For this latter reason such a motion of a point, periodic in respect to time, is often called simple harmonic, and the mathematical function which describes this motion is called a circular function.

Simple Harmonic Motion. Consider Fig. 52, and imagine a point P moving with a uniform speed round the circumference of a circle of radius a , and having centre O , and suppose that the point executes n complete revolutions per second. This means that reckoning time from the instant that the moving point P coincides with A on the reference line OX , the point will next occupy this position after the lapse of $\frac{1}{n}$ seconds. During this time the line OP , joining P to the centre of the circle, will have swept through four right-angles or 2π radians of angle. As the circular motion of P is

uniform we can say that in a time interval $\frac{1}{2n\pi}$ the line OP will arrive from OA to a position OP_1 such that the geometrical angle AOP_1 is 1 radian. The quantity $2\pi n$ which represents the number of radians of angle swept out by the moving line OP in one second is called the angular velocity of the rotation of P and this quantity is usually denoted by the Greek letter ω . As the execution of 1 radian of angle is the result of P moving a distance of a along the circumference, it follows that the velocity of P in length units per second is $a\omega$.

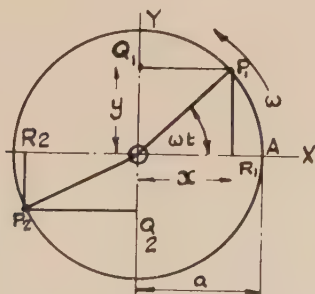


FIG. 52.

Let us think first of a position of P like P₁ where OP₁ makes a geometrical angle, P₁OA with OX. After the execution of exactly 1 revolution, or after the lapse of a time $\frac{1}{n}$

seconds, the point will be at the same position. $\frac{1}{n}$ can,

therefore, be called the periodic time of rotation, and denoted by T . The geometrical angle $\angle P_1OP$ will be the same at the time $t + T$ as it was at T , but during T the angle swept out by OP will have increased by 2π radians. This is a wider conception of angle than that of the mere measure of the inclination of two lines, and this conception, that of an amount of rotation, measured by the number of radians swept through by a line like OP_1 , is very important. According to this conception the angle between OP and the reference position OA

is continually increasing as P revolves uniformly, the increase being at the rate of ω or $2\pi n$ radians a second. Angles described by rotation in this way are reckoned positive when the rotation is anticlockwise as shown, and negative when the rotation is clockwise.

Let Q be a point on the vertical line OY, located, by drawing from P, a line parallel to OX. As P revolves uniformly the point Q will evidently oscillate about O in the vertical diameter of its travel. The distance y of Q from O, called the projection of OP on OY, and in accordance with a convention previously defined, reckoned positive above and negative below OX, will continuously vary between the limits $+a$ and $-a$. The motion of Q is periodic, and its amplitude is a . Let us think again of a time instant, t , which makes P take the position P_1 . The distance y is plainly $a \times \sin P_1OA$, for the sine of this geometrical angle is $\frac{OQ_1}{a}$. We define the sine

ratio as the value $\frac{y}{a}$ for all angles, in the extended sense of

measures of radian rotation. Thus the ratio $\frac{OQ_1}{a}$ is the sine not only of the geometrical angle P_1OA , measured in radians by ωt , but of all angles equal to ωt plus any exact number of revolutions, or of $(\omega t + m \times 2\pi)$ where m is any whole number. Thus, we may say $\sin \omega t = \sin (\omega t + m \times 2\pi)$, or as 2π , one revolution, is swept out in T seconds, so that it is equal to ωT , $\sin \omega t = \sin \omega(t + mT)$. Thus $\sin \omega t$ agrees with our basic definition of a periodic function; the periodic movement of Q, the projection of P on the vertical OX, can be symbolically stated $y = a \sin \omega t$ and the position of Q, and the distance y , are exactly reproduced when the time t is increased by any exact multiple of T. The motion of the point Q is said to be simple harmonic, its amplitude or maximum displacement from the mean position is a , its frequency is n , and its periodic time $T = \frac{1}{n}$, and the manner of its variation, depending upon the sine of a varying angle is said to be sinusoidal. The geometrical angle POA, which is the remainder when ωt is divided by 2π , is called the phase of y .

Circular Functions. In Chapter III the sine of an angle is defined as a ratio that fixes the shape of a right-angled triangle. According to the extended idea of angle, $\sin \theta$ is a periodic function of θ , an amount of turning which may be reckoned in revolutions, degrees, or radians, and which can be conceived to increase indefinitely. The function $y = a \sin \theta$ or $a \sin \omega t$ can be represented graphically by the obvious geometrical construction of Fig. 53, in which a portion of the horizontal axis of time or angle, representing T or 2π , is divided into as many equal parts as is the circumference of a circle of a radius representing the amplitude a . This graph shows how $\sin \theta$ changes sign as θ is continually increased. These changes are often expressed in the following way. The quadrants of the circle of Fig. 52 may be numbered 1 to 4, starting from OA and reckoning in a positive direction. Thus OP_1 is in the first quadrant and OP_2 in the third. With this convention we can say that when the residual angle obtained by dividing θ by 2π lies in the first and second quadrants, $\sin \theta$ is positive. As negative rotation of the line OP changes the direction of movement of Q , we have $\sin (-\theta) = -\sin \theta$.

If in Fig. 52 the point R is located by a perpendicular from P to the horizontal OX , it is evident that the movement of R will be periodic. The distance of x of R from O positive when R is to the right and negative when it is to the left of O is evidently $a \cos POA$, and the movement of R can be represented by the functions $x = a \cos \omega t$. This function can be represented in graphical form in a manner similar to that employed for $a \sin pt$. The graph of $y = a \cos \omega t$ or $a \cos \theta$ is shown also in Fig. 53. At a time O when P coincides with A , $a \sin \omega t = 0$, but x or $a \cos \omega t$ has its maximum value of a . Also when $a \sin \omega t$ has its maximum and P is on the vertical line x_1 or $a \cos \omega t$ is zero. It is evident that when P lies to the left of the vertical OY , or in the second and third quadrants x , and hence $\cos \omega t$ is negative. It is also evident that the position of R , and the sign of x is the same whether rotation of OP is positive or negative, hence $\cos -\theta = \cos \theta$. Further, and what is most important, the cosine of an angle θ is the same as the sine of an angle, $\frac{1}{4}$ turn, or $\frac{\pi}{2}$ radians, or 90 degrees

greater than θ . Symbolically $\cos \omega t = \cos \left(\omega t + \frac{\pi}{2} \right)$. Similarly $\sin \omega t = \cos \left(\omega t - \frac{\pi}{2} \right)$.

The graph of $\tan \theta$ can be traced by the geometrical construction shown in Fig. 54. Here at the end A of the diameter of a circle of radius a , a perpendicular line AB is drawn, and the circumference of the circle is equally divided as in

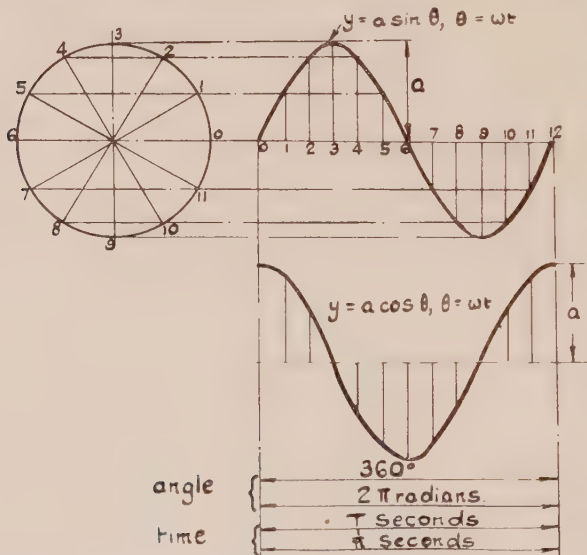


FIG. 53.

Fig. 53. The tangent of the angle of any of the radii determined by the points of division is equal to the length like AS, where S is the point where the radius meets the perpendicular line. It is evident that, as the moving point P of Fig. 52 approaches the vertical line OY, the prolongation of OP will become nearer and nearer to parallelism with a vertical line like AB of Fig. 54. Thus, as θ approaches $\frac{\pi}{2}$ radians or 90 degrees, $\frac{3\pi}{2}$ radians, 270 degrees, or $\frac{3}{4}$ turn, $\tan \theta$

increases without limit. We note that if the moving point P of Fig. 54 goes to the 3rd quadrant both the base and vertical side of the right-angled triangle determining the tangent ratio are negative, and the tangent is the ratio of two negative quantities and is positive; $\tan \theta$ is thus positive when the remainder obtained by dividing θ by 2π lies in the first or third quadrants. The graph shows clearly two features in which the $\tan \theta$ function differs from $\sin \theta$ and $\cos \theta$. First its

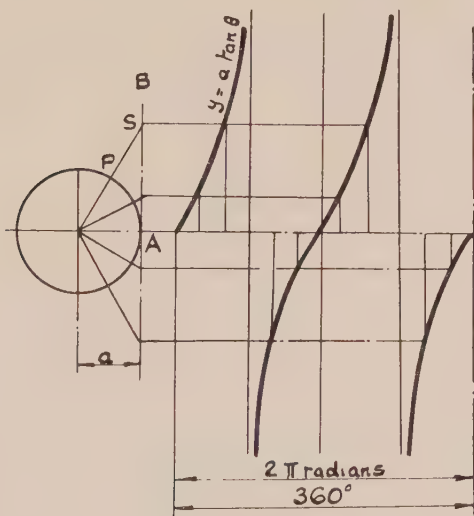


FIG. 54.

value is exactly repeated when θ is increased by any multiple of π , that is, by any exact number of half-turns of angle, so that its period is one half that of $\sin \theta$ and $\cos \theta$. Secondly, $\tan \theta$ is a function discontinuous when θ has the value of $\frac{\pi}{2} + m\pi$, where m is any whole number.

We may now summarise the information we have obtained about the circular functions $\sin \theta$, $\cos \theta$, and $\tan \theta$. They are all periodic, in regard to the independent variable θ , the period for $\sin \theta$ and $\cos \theta$ is 2π , and for $\tan \theta$ is π . $\sin \theta$ and $\cos \theta$ are functions everywhere continuous. Since

$\cos(-\theta) = \cos \theta$ the cosine is an even function, and as $\sin(-\theta) = -\sin \theta$ the sine is an odd function. The sine and cosine functions have a maximum numerical value of 1, the numerical value of the tangent function is unlimited. It follows from the last statement that every positive number is the tangent of some angle lying between the values 0 and $\frac{\pi}{2}$ radians (90 degrees).

Special Values of the Circular Functions. If the reader remembers the metrical properties of special right-angled triangles having acute angles of 45 and 60 degrees, like ordinary set-squares, which were discussed in Chapter III,



FIG. 55.

he will see at once by considering Fig. 55 that we can write down values of the circular functions for angles of $\frac{\pi}{6}$, $\frac{\pi}{4}$, and $\frac{\pi}{3}$ radians, or 30, 45, and 60 degrees. These values are included in the following table, and the reader should not only see that they are correct, but he should be able to recall them immediately, preferably by a mental picture of the corresponding triangles of Fig. 55.

θ		$\sin \theta$	$\cos \theta$	$\tan \theta$
Radians	Degrees			
0	0	0	1	0
$\pi/6$	30	$\frac{1}{2}$	$\sqrt{3}/2$	$1/\sqrt{3}$
$\pi/4$	45	$1/\sqrt{2}$	$1/\sqrt{2}$	1
$\pi/3$	60	$\sqrt{3}/2$	$\frac{1}{2}$	$\sqrt{3}$
$\pi/2$	90	1	0	Infinite
π	180	0	-1	0
$3\pi/2$	270	-1	0	Infinite
2π	360	0	1	0

Consider Fig. 56 where $\angle BOA = \theta$ is a small angle defined by the arc AB and the radius OA or OB which is equal to 1, in some length unit. The length of the arc BA measured in terms of this length unit is the radian measure of θ , and the length of the perpendicular BS or OA is the sine of θ . We see that, as θ gets smaller and smaller, the arc BA and the sine PR become more and more nearly equal in length, so that the smaller θ is, the more nearly $\sin \theta = \theta$. Expressed more exactly, the limiting value of the ratio $\frac{\sin \theta}{\theta}$ as θ approaches zero is equal to 1. Thus when an angle θ is small its sine is approximately equal to its radian measure, and as π or 3.14159



FIG. 56.

radians make up 180 degrees, so that 1 radian is very nearly 57.3 degrees, the sine of a small angle is about equal to its degree measure divided by 57.3.

We have seen on p. 45 that $\cos^2 \theta = 1 - \sin^2 \theta$, where $\cos^2 \theta$ and $\sin^2 \theta$ stand respectively for the squares of the sine and cosine. For small angles, therefore, $\cos^2 \theta = 1 - \theta^2$. Now the

square of $\left(1 - \frac{\theta^2}{2}\right)$ is equal to $1 - \theta^2 + \frac{\theta^4}{4}$, and when θ is small

θ^4 will be negligible in relation to θ^2 , so that $(\cos^2 \theta)^2 = \left(1 - \frac{\theta^2}{2}\right)^2$

very approximately. A fair approximation for $\cos \theta$ when θ

is small is 1, a better approximation is $1 - \frac{\theta^2}{2}$. The reader

should note and thoroughly understand that the approximations for small angles $\sin \theta = \theta$ and $\cos \theta = 1 - \frac{1}{2}\theta^2$ apply only when θ is measured in radians.

Compounding Perpendicular Simple Harmonic Motions.

Referring back to Fig. 52 we see that the two simple harmonic motions $x = a \cos \omega t$, and $y = a \sin \omega t$, define two simultaneously

varying directed steps which in the complex number notation may be denoted by $a \times \cos \omega t$ and $a \times j \sin \omega t$, and that at any instant these two steps make up a radial step representing a complex number of amplitude a and argument or angular inclination equal at this instant to ωt . If a is 1 then, considering variations of the angle $\theta = \omega t$ we have z , the complex number, is equal to $\cos \theta + j \sin \theta$, as we have seen before. We see now that considering θ as an independent variable, z is a function of θ , say $f(\theta)$, and it has the property that $f(\theta_1) \times f(\theta_2) = f(\theta_1 + \theta_2)$, for, according to the convention for the multiplication of complex numbers, a unit step inclined θ_1 multiplied by a unit step inclined θ_2 is a unit step inclined $(\theta_1 + \theta_2)$. $f(\theta)$ is therefore something like the function $A^x = \psi(x)$ for as $A^{x_1} \times A^{x_2} = A^{x_1 + x_2}$, $\psi(x_1) + \psi(x_2) = \psi(x_1 + x_2)$, so that the θ of $\cos \theta + j \sin \theta$ has one of the properties of the power index in A^x . The θ is, however, inherently different from an ordinary real index of a power, for, first, 1^x is always equal to 1, and secondly, A^x , when A is not equal to 1, is always changed in magnitude when x is changed. A change of θ does not alter the magnitude of the complex number z . Now consider Fig. 57, OA represents a real number, and raising



FIG. 57.

this to a small positive power will slightly increase it to a value like OA_1 , the alteration geometrically is in the same direction as OA. Now think of a number $f(\theta)$ when $\theta = 0$, this will be represented also by a line like OA in the real number axis. If θ takes a very small value, then $f(\theta)$ will be represented by the line OR equal to, but inclined relatively to, OA by this small angle. The change has been brought about by adding to the step OA a very small step, not like AA_1 in line with it but AR at right-angles to it. We might consider, therefore, that θ has the property of a power index, the increase of which continually changes Z in a direction at right-angles to itself, and this might be expressed by saying that $Z = A^{j\theta}$. This will work, for $A^{j\theta_1} \times A^{j\theta_2} = A^{j(\theta_1 + \theta_2)}$ by the index

law of algebra, but to say, at present, that $\cos \theta + j \sin \theta = A^{j\theta}$ is not intelligible, because the nature of the number A must be elucidated. We shall return to this matter in a future chapter.

If we think of two simultaneous simple harmonic motions of equal amplitude in the purely geometrical sense, that one $x = a \cos \omega t$ defines the x co-ordinate, and the other $y = a \sin \omega t$ the y co-ordinate of a graph, then combining these equations we have at once $x^2 + y^2 = a^2(\cos^2 \omega t + \sin^2 \omega t) = a^2$, the equation to a circle discussed in the early part of Chapter V. We may note, however, that the ordinary equation to a circle $x^2 + y^2 = a^2$ is the result of combining two separate equations $x = a \cos \theta$ and $y = a \sin \theta$; these latter equations are called the parametric equations of the circle; they give the functional relationship of x and y , not directly but in terms of their relations to another variable or parameter θ .

If the perpendicular simple harmonic motions are of different amplitude, then they will define the co-ordinates of a curve the parametric equations of which are $x = a \cos \omega t$ and $y = b \sin \omega t$ where b is the amplitude of the vertical motion.

These equations are the same as $\frac{x}{a} = \cos \omega t$, $\frac{y}{b} = \sin \omega t$, whence

by squaring and combining them we obtain the resulting equation $\frac{x^2}{a^2} + \frac{y^2}{b^2} = \cos^2 \omega t + \sin^2 \omega t = 1$ which, as we saw on

p. 128, is the equation to an ellipse. This is shown in Fig. 58 where the motion $x = a \cos \omega t$ is determined by the circular revolution of P_1 and $y = b \sin \omega t$ by that of P_2 , and P is the point in the ellipse corresponding to the position OP_2P_1 of the line rotating with angular velocity ω . The line joining the centre of the circle to the point P is called a radius vector of the ellipse; it evidently rotates continuously with the same frequency as OP_2P_1 but not uniformly, although it coincides in direction with OP_2P_1 in the horizontal and vertical directions. The angle ϕ between OP and OA has for its tangent

the ratio $\frac{y}{x}$ or $\frac{b \sin \omega t}{a \cos \omega t}$ which, as shown on p. 46, is equal to

$\frac{b}{a} \tan \omega t$. The angular motion of the variable line OP is

therefore defined by the equation $\tan \phi = \frac{b}{a} \tan \omega t$, the angular velocity of this line is variable, being sometimes less and sometimes greater than the constant angular velocity ω of

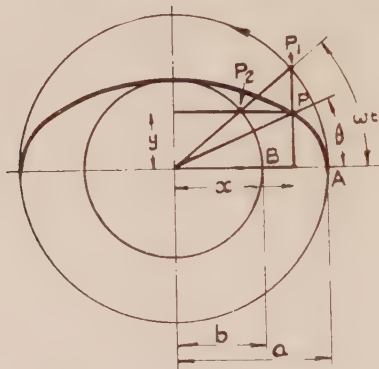


FIG. 58.

OP_2P_1 , but the average angular velocity of OP over a whole rotation is equal to ω .

Compounding Simple Harmonic Motions in the same Directions. The instantaneous sum of two simple harmonically varying distances $x_1 = a_1 \sin \omega t$ and $x_2 = a_2 \sin \omega t$ which pass through zero values synchronously, and of the same frequency, is $x = (a_1 + a_2) \sin \omega t$, a synchronous simple harmonically varying distance of amplitude equal to the sum of the component amplitudes. Suppose, however, that the two periodic functions such that one attains zero value after the other. This state of affairs is represented graphically in Fig. 59. Here the circular rotation of a line OP_1 defines a simple harmonic displacement $y_1 = a \sin \omega t$, and the rotation of another line OP_2 defines a second simple harmonic displacement y_2 . If the angle between OP_1 and OP_2 is ϕ , the angle ωt of OP_1 will always be greater by ϕ than that of OP_2 , so that $y_2 = b \sin (\omega t - \phi)$, which means that y_2 attains zero values by a time interval t_1 , after y_1 which is given by the relation $\frac{\phi}{2\pi} = \frac{t_1}{T}$,

T being the period time. A time interval of this kind is called a difference of phase, because the phase of y_2 is always ϕ less than that of y_1 , and as the variations of y_2 are later than those of y_1 , y_2 is said to lag in phase on y_1 while y_1 leads in phase on y_2 . It is easy to see from the diagram that at any instant of time the sum of the y_1 and y_2 distances is equal to the simple harmonically varying distance y determined by a rotating line which is the diagonal of a parallelogram whereof OP_1 and OP_2 are adjacent sides, for as

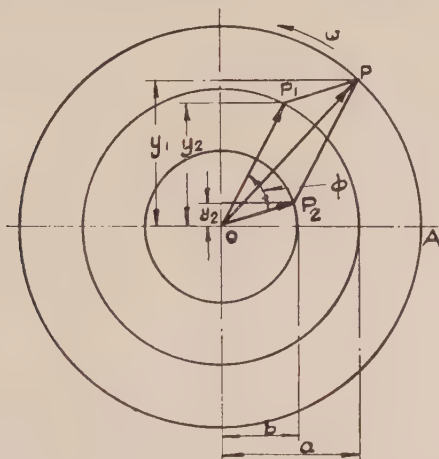


FIG. 59.

P_1P is parallel to OP_2 the excess of y over y_1 is manifestly equal to y_2 . Further, this is true, whatever may be the position of OP_1 and OP_2 , for as they rotate synchronously the diagonal OP of the parallelogram will rotate with them and this parallelogram will keep its shape unchanged. Of course if y_1 is negative while y_2 is positive these distances will be added or combined algebraically and their sum or resultant will be an arithmetical difference. We thus arrive at the important rule that the sum of two simple harmonic periodic functions of the same frequency of amplitudes a and b , having a phase difference of angular measure ϕ , is given in magnitude and phase by the diagonal of a parallelogram having adjacent

sides a and b at an angle of ϕ , and passing through this angle. Thus simple harmonic periodic functions can be added or combined as directed steps or vectors. This is shown in Fig. 60. Here OA is a line drawn to scale to represent the amplitude a , and OB is drawn at an angle of ϕ to OA , in a negative direction to represent a lag in phase, the length OB representing the second amplitude b to scale. The resultant simple harmonic function is represented by the vector OP , which indicates that this resultant has an amplitude given to scale by the length OP , and that the phase of the resultant is ϕ_1 lagging on the component represented by the vector OA . Instead of drawing a parallelogram we could have represented the vector OB by a line AP parallel and equal to it and so have

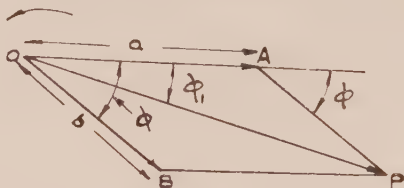


FIG. 60.

obtained the resultant vector OP by constructing the triangle OAB only. The difference between two simple harmonic functions is obtained by a rule similar to that explained on p. 77 for obtaining the difference of two directed steps. The vector of the function to be subtracted and the reversed vector is added.

If the two simple harmonic periodic functions are not of the same frequency the nature of their sum is not so clearly elucidated. We can, however, see in a general way what this will be, by the vector construction which has been deduced for we can consider that the phase of the function with the greater frequency is continually advancing in a leading direction. Consider Fig. 61 and let the vector OA represent the amplitude of the function with the smaller frequency. If a circle of radius equal to the amplitude of the second function be described with A as centre, then this function can be represented by a radius such as AP_2 which, rotating at a

speed corresponding to the difference of the two frequencies, gives a resultant OP , the amplitude of which continually varies between the limits of the sum and difference of the component amplitudes with a frequency corresponding to this difference. This is the explanation of the phenomenon of beats observed when two musical tones of nearly the same pitch produce a resultant tone of continually waxing and waning loudness.

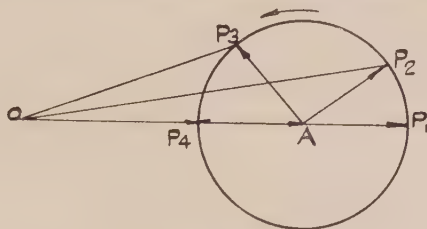


FIG. 61.

Circular Functions of Compound Angles. It has been shown on p. 82 that it follows from the multiplication rule for complex numbers that as

$$\cos (\theta_1 + \theta_2) + j \sin (\theta_1 + \theta_2) = (\cos \theta_1 + j \sin \theta_1)(\cos \theta_2 + j \sin \theta_2)$$

$$\text{and } \cos (\theta_1 - \theta_2) + j \sin (\theta_1 - \theta_2) = (\cos \theta_1 + j \sin \theta_1)(\cos \theta_2 - j \sin \theta_2)$$

we can obtain the following formulae :

$$\cos (\theta_1 \pm \theta_2) = \cos \theta_1 \cos \theta_2 \mp \sin \theta_1 \sin \theta_2$$

$$\text{and } \sin (\theta_1 \pm \theta_2) = \sin \theta_1 \cos \theta_2 \pm \cos \theta_1 \sin \theta_2.$$

where the double signs \pm and \mp in the first equation means that a $+$ sign on the left-hand side is taken with a $-$ sign on the right, and vice versa. These formulae are proved geometrically in the textbooks on trigonometry which also contain and prove a large number of other formulae for circular functions of two angles. All these formulae or identities are very important for those who desire to attain to skill in mathematical technique, but they are difficult to remember unless continually used. Those given above should, however, be memorised by all who wish to understand mathematical symbolism, and the following important derived identities

which are particular cases of $(\theta_1 + \theta_2)$ when the two angles are equal

$$\sin 2\theta = 2 \sin \theta \cos \theta$$

and

$$\cos 2\theta = \cos^2 \theta - \sin^2 \theta$$

As $\cos^2 \theta + \sin^2 \theta = 1$ for all values of θ , this last equation can be converted to the following important formulae :

$$\sin^2 \theta = \frac{1}{2}(1 - \cos 2\theta)$$

$$\cos^2 \theta = \frac{1}{2}(1 + \cos 2\theta).$$

These last equations are of particular interest for, if θ is a variable angle, ωt , they show that $\sin^2 \omega t = \frac{1}{2} - \frac{1}{2} \cos 2\omega t$, that is to say, the square of a simple harmonic function is

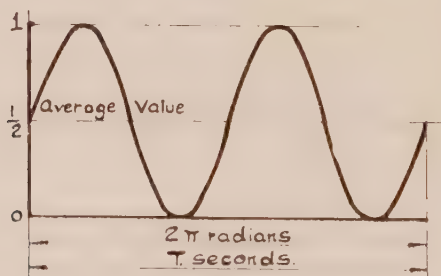


FIG. 62.

$\sin^2 \theta$, starting at $\theta = 45^\circ$ or $\frac{1}{4}\pi$ radians, is given in Fig. 62. This shows the constant and double frequency components. The graph lies entirely above the horizontal axis of co-ordinates, because as the square of a negative quantity is positive, $\sin^2 \theta$ always has a positive value.

Average Values of Periodic Functions. The conception of the average value of a periodic function taken over a complete exact cycle of change is very important in applied science and engineering. For example, the heating effect of an alternating current must be continually varying as the current itself is very approximately a simple harmonic periodic function of the time, having, in this country, a frequency of 50 per second, but to the user of electricity supply it is the average heating

effect of the current that gives its utility for cooking, lighting, or practical heating purposes.

It is easy to see that the average value of a simple harmonic periodic function $a \sin \omega t$ taken over a complete cycle is zero, for the change in the function from 0 through $+a$ to 0 in the first half cycle is exactly repeated in the opposite sense from 0 through $-a$ to 0 in the second half cycle. In short, to every positive value of the function in the first half cycle there corresponds an equal negative value in the next, so that throughout the whole cycle positive and negative values cancel.

It is not difficult to deduce the average value of the square of a simple harmonic function. For as $\sin^2 \omega t + \cos^2 \omega t = 1$ for all values of the angle ωt , and as $\cos \omega t = \sin \left(\omega t - \frac{\pi}{2} \right)$ we have $a \sin^2 \omega t + a \sin^2 \left(\omega t - \frac{\pi}{2} \right) = a$. Now, as $a \sin^2 \omega t$ and $a \sin^2 \left(\omega t - \frac{\pi}{2} \right)$ are the squares of similar functions differing only in phase, their average values must be equal, and as their sum is always a , the sum of these equal average values must be a , so that the average value of $a \sin^2 \omega t = \frac{a}{2}$, and this is the average value of $a \cos^2 \omega t$. The foregoing reasoning is confirmed by using the formula for $\sin^2 \theta$ in the last section, for as $a \sin^2 \theta = \frac{a}{2} - \frac{a}{2} \cos 2\theta$, and as the average value of the double frequency term $\cos 2\theta$ over a complete cycle of change of $\sin \theta$ is zero, $a \sin^2 \theta$ has an average value equal to the constant quantity $\frac{a}{2}$.

If, however, we consider the average value not of the product of a sine by a sine function, but of $\sin \omega t$ and $\cos \omega t$, we obtain a different result, for, using the formula for $\sin 2\theta$ in the last section we see that $\sin \omega t \times \cos \omega t = \frac{1}{2} \sin 2\omega t$ which is a simple harmonic function of double frequency. The average value of $\sin \omega t \times \cos \omega t$ is therefore zero.

Now consider the product $a \sin \omega t \times \sin (\omega t - \phi)$, that of two simple harmonic functions differing in phase by ϕ and of the same frequency. We can easily discover the average

value of this product by using a formula of the last section to change $\sin(\omega t - \phi)$, for with this change we find that the product is

$$a \sin \omega t (\sin \omega t \cos \phi - \cos \omega t \sin \phi)$$

or

$$a \sin^2 \omega t \cos \phi - a \sin \omega t \cos \omega t \sin \phi.$$

Now ϕ is a constant angle representing phase difference, so that $\cos \phi$ and $\sin \phi$, like a , are both constant. The average value of the first term in the last expression is $a \cos \phi$ multiplied by the average value of $\sin^2 \omega t$. This is $\frac{1}{2}a \cos \phi$. The average value of the second term depends on the average value of $\sin \omega t \times \cos \omega t$ which we have found to be zero, so the average value of the whole term is zero. The average value of the product $a \sin \omega t \sin(\omega t - \phi)$ is therefore $\frac{1}{2}a \cos \phi$.

As ϕ varies from 0 to $\frac{\pi}{2}$, the average value of the product varies from that of $a \sin \omega t \sin \omega t$ or $a \sin^2 \omega t = \frac{a}{2}$ which, of course, is $\frac{1}{2}a \cos 0$, to the average value of $a \sin \omega t \sin \left(\omega t - \frac{\pi}{2} \right)$ or $a \sin \omega t \cos \omega t$ which is zero, and which is equal to $\frac{1}{2}a \cos \frac{\pi}{2}$.

The geometrical meaning of this is that the average line of Fig. 62, which corresponds to $\sin \omega t \sin(\omega t - \phi)$, when $\phi = 0$, descends as ϕ in $\sin \omega t \sin(\omega t - \phi)$ increases numerically, and that the height of this line above the horizontal axis is equal to $\frac{1}{2} \cos \phi$. When $\phi = \frac{\pi}{2}$ and $\cos \phi = 1$ the line of symmetry

becomes coincident with the horizontal axis of co-ordinates and the function $\sin \omega t \sin(\omega t - \phi)$ becomes a simple harmonic function of double frequency with an average value of zero.

Now let us return to a consideration of the formula for $\cos(\theta_1 - \theta_2)$ given in the last section, and let us consider that θ_1 and θ_2 are angles continually varying in time, but with different angular velocities ω_1 and ω_2 , $\cos(\theta_1 - \theta_2)$ can then be written $\cos(\omega_1 - \omega_2)t$, and its equivalent expanded form becomes $\cos \omega_1 t \cos \omega_2 t + \sin \omega_1 t \sin \omega_2 t$. $\cos(\omega_1 - \omega_2)t$ is a simple harmonic function, having a frequency equal to the difference of the frequencies corresponding to ω_1 and ω_2 , and we know of this function that, provided ω_1 is not equal to

ω_2 , however small the resultant frequency will be, that its average value taken over a whole cycle of its change must be zero. If, however, the two component frequencies are equal then $\cos (\omega_1 - \omega_2)t = \cos 0 = 1$, as $\omega_1 = \omega_2$. Now let us consider the equivalent form of this function, $\cos \omega_1 t \cos \omega_2 t + \sin \omega_1 t \sin \omega_2 t$. We can infer that the average values of these two products will be the same, for the cosine functions of both $\omega_1 t$ and $\omega_2 t$ are simply sine functions, $\frac{\pi}{2}$ later in phase.

It thus follows that provided ω_1 and ω_2 are not equal then the average values of $\sin \omega_1 t \times \sin \omega_2 t$ and $\cos \omega_1 t \times \cos \omega_2 t$ taken over a complete cycle of change corresponding to the difference of the two frequencies corresponding to ω_1 and ω_2 are both zero, as the average value of $\cos (\omega_1 - \omega_2)t$ is zero. If, however, $\omega_1 = \omega_2$, then as $\cos (\omega_1 - \omega_2)t$ becomes $\cos 0 = 1$ the equal average values of $\sin \omega t \times \sin \omega t = \sin^2 \omega t$ and $\cos^2 \omega t$ total 1, and each of these averages is $\frac{1}{2}$, as we found before. The average value of the product of two simple harmonic functions of different frequencies is therefore zero. This result is very important in electrical technology; it shows that, as the torque of an electric motor depends upon the product of the current in its windings and its magnetic field, if current and field are alternating at different frequencies the machine can produce no average torque.

Velocity and Acceleration of Simple Harmonic Motion. Let us return to our original definition of simple harmonic motion, as the movement of a point Q (Fig. 68), which is the projection on the vertical of a point P revolving in a circular path of radius a with a constant angular velocity of ω radians per second. Q executes an oscillatory motion in the vertical diameter of the circle; its speed evidently varies; at the ends of its travel the direction of its movement reverses and it must be momentarily at rest. In the middle of its travel Q is moving parallel to P and its speed here would seem to be the greatest. Let us consider an instant defined by time t , which defines the position of P by the angle $POA = \omega t$. The linear speed of P is constant, that is, if a is in centimetres it travels through ωa centimetres every second. Suppose that, at the instant t , P continued to move, not in its circular path but in

the direction it was travelling at this instant, it would go from P to P_1 in one second where the line PP_1 is tangential to the circle, at right-angles to the radius, and the length PP_1 is equal to the linear speed ωa . Thus tangential motion would take the point Q to Q_1 , which is equal to $\omega a \cos \omega t$, as PP_1 is inclined at an angle ωt to the vertical diameter. This distance $\omega a \cos \omega t$, which would have been traversed by the point Q

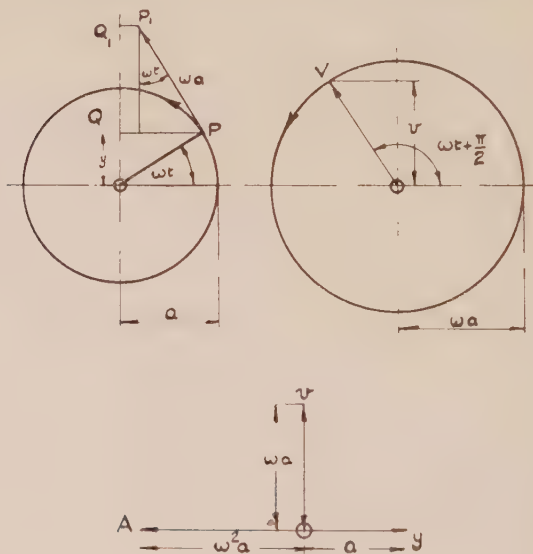


FIG. 63.

in one second from the instant t had this movement of P continued unchanged in direction along the tangential line, we shall call the velocity of Q at the instant t , and we see from the right-hand diagram that as the linear speed of P can at any instant be represented by a rotating line OV of radius ωa , $\frac{1}{4}$ revolution or $\frac{\pi}{2}$ in angle in advance of OP , this velocity v of Q can be expressed generally as $\omega a \cos \omega t$, for v is the projection of OV on the vertical, the angular position of OV is $\omega t + \frac{\pi}{2}$, and $\sin \left(\omega t + \frac{\pi}{2} \right)$ is equal to $\cos \omega t$. Thus, using the

vector method of representing simple harmonic periodic function, if Oy in the lower diagram is the vector of the variable displacement of Q having a length a , then Ov a vector $\frac{\pi}{2}$ radians or 90 degrees leading it in phase and of length ωa , is the vector of the velocity of Q .

If T is the periodic time of the simple harmonic motion, then as 2π radians are swept out by OP in T seconds $\omega = \frac{2\pi}{T}$, so that the velocity of Q at any instant is $\frac{2\pi a}{T} \cos \omega t$. 2π and $\cos \omega t$ are both ratios or pure numbers dimensionally so that the formula gives a velocity of the dimensions of a length a divided by a time T , as it ought, according to the explanation of dimensions on p. 88.

The maximum velocity of Q is at the instant that $\cos \omega t$ is equal to 1; it is then ωa , or $\frac{2\pi a}{T}$. Q travels from one end of a diameter of the path of P to the other, a distance $2a$, in half the periodic time $\frac{T}{2}$ seconds; its average speed is therefore $\frac{4a}{T}$. The ratio of the average to the maximum velocity is thus $\frac{4a}{T} \div \frac{2\pi a}{T} = \frac{2}{\pi}$.

Let us consider this idea of the speed of the point Q executing simple harmonic motion from another point of view. Fig. 64 is a part of a graph representing the movement of Q as time varies. Any instant of time t defines a distance y . If t increases by an amount δt , y will change by a corresponding amount δy , and the average velocity of the point Q in the interval is distance moved divided by the time interval in which this movement took place. The average speed is proportional to the slope of the line AB joining points on the graph determined by $t = (t + \delta t)$, and the slope of the graph as defined on p. 120 at the instant t is the limiting position which the AB tends to assume as the time interval is decreased. The straight line with this slope will be the tangent at the point defined by t . If the slope of this graph is some number

n , then this means that the tangent rises a distance n units on the graph diagram for one unit horizontal travel. If one length unit on the graph represents l units of travel of Q , and one unit also represents s seconds of time, then the speed of the simple harmonic motion is the slope $n \times \frac{l}{s}$.

We can work out this idea of simple harmonic velocity in yet another way. At time t the displacement Q is $a \sin \omega t$. If t increases by a small amount δt the displacement will become $a \sin \omega(t + \delta t) = a \sin (\omega t + \omega \delta t)$. We can expand this sine of

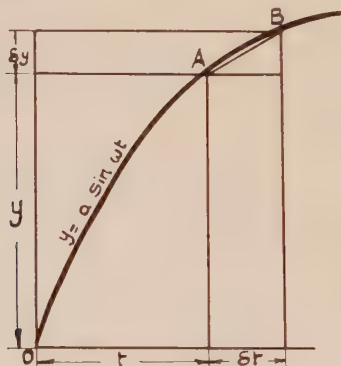


FIG. 64.

the sum of two angles by the formula of p. 171 and the new displacement becomes $a \sin \omega t \cos (\omega \delta t) + a \cos \omega t \sin (\omega \delta t)$. Now $\omega \delta t$ is a small angle, and by the reasoning of p. 165 we have seen that the smaller it becomes the closer does $\cos (\omega \delta t)$ approach 1 and the closer does $\sin (\omega \delta t)$ approach the radian measure $\omega \delta t$ of the angle. Thus the smaller the increase δt of time becomes, the nearer does the new displacement approach the value $a \sin \omega t + a \cos \omega t \times \omega \delta t$, and the nearer does the change of displacement in the small interval δt approach the value $\cos \omega t \times a \omega \delta t$. The ratio, change of displacement divided by time, is the average velocity in the interval δt , and it follows therefore that as this interval is diminished the average velocity in the interval tends to the limiting value, $a \omega \cos \omega t$. This limiting value is called the velocity at time t .

We arrive therefore at the important conclusion that the velocity of a point Q executing periodic simple harmonic motion is also simple harmonic. The phase of the velocity, which may be called rate of change of position, is $\frac{\pi}{2}$ radians or 90 degrees leading on the phase of displacement. Thus velocity with simple harmonic motion is continually changing. Now rate of change of velocity is called acceleration, and we can easily work out the value of this acceleration. For if the rate of change of position $a \times \sin \omega t$ is $a\omega \times \sin \left(\omega t + \frac{\pi}{2} \right)$ we can infer that the rate of change of $a\omega \times \sin \left(\omega t + \frac{\pi}{2} \right)$ will be $a\omega \times \omega \times \sin \left(\omega t + \frac{\pi}{2} + \frac{\pi}{2} \right)$, a simple harmonic function with amplitude ω times that of velocity and $\frac{\pi}{2}$ in phase advance of it. The vector of acceleration is therefore OA in the lower diagram of Fig. 63, having a length $\omega^2 a$ and being π radians, or half a turn in phase advance of the displacement vector Oy. This shows that the acceleration with simple harmonic motions varies in direct proportion to the displacement but is opposite in sign to it, and we deduce this from the formula, for $a\omega^2 \sin \left(\omega t + \frac{\pi}{2} + \frac{\pi}{2} \right) = \omega^2 a \sin (\omega t + \pi)$, and it is easy to see that $\sin (\omega t + \pi)$ is always the negative of $\sin \omega t$. Thus, as the displacement $y = a \sin \omega t$, the acceleration is $-\omega^2 y$, the negative sign showing that when y is positive the velocity in this direction is decreasing.

We can easily interpret this dynamically, for if a body executing a simple harmonic motion has a mass m , a central force will be required to restrain its motion and to keep it oscillatory, and this force is equal to $m \times \text{acceleration}$, by Newton's second law. The measure of the central force is, therefore $m\omega^2 y$. It is proportional to the displacement y . The converse of this statement is true, that if a body with mass is constrained by force proportional to its displacement towards a central position, always to return to this position, then the body, if displaced and released, will oscillate with

simple harmonic motion. The control force per unit displacement being constant, can be written $f = \frac{\text{force}}{y}$, and as $\text{force} = m\omega^2 y$ we have constant $f = m\omega^2$, and since $\omega^2 = 4\pi^2 n^2$ the frequency of the oscillation n is equal to $\frac{1}{2\pi} \sqrt{\frac{f}{m}}$. Thus knowing the mass of a body and its elastic constraint expressed in force per unit displacement, the frequency of its oscillation can be calculated.

Fourier's Theorem. We shall conclude this chapter by a brief reference to a very far-reaching mathematical theorem, due to Fourier, which states that any periodic function defined by $f(t) = f(t + mT)$ where m is a whole number is equal to the sum of simple harmonic or circular periodic functions, the frequencies of which are whole-number multiples of the frequency $\frac{1}{T}$ of the principal function. Thus any periodic function can be expressed as the sum of a number of terms like $a \sin(p\omega t + \theta)$ where θ is the phase difference of the simple harmonic component and the principal function and p is a whole number. The frequency of this component is p times that of the principal function. These component circular functions are known as harmonics, and the number p gives the order of the harmonic. When $p=1$ the frequency of the harmonic is equal to that of the principal function, this is called the first harmonic or, more generally, the fundamental. When $p=2$ the corresponding component is called the second harmonic, and so on.

It is not difficult to see by a simple geometrical construction how non-sinusoidal periodic functions can be made up of simple harmonic components. In Fig. 65, for instance, the dotted graphs are those of a fundamental and a third harmonic, and the addition of ordinates gives the full-line graph, which has a flat topped and dimpled geometrical wave form. If the third harmonic component is drawn so that at the start of the graph it is increasing in the negative instead of the positive direction, then it is not difficult to see that the peak of the fundamental will be reinforced by a peak of the third harmonic, and the geometrical wave form will be peaked. In

each of these cases, however, the half-wave is symmetrical in shape, and successive half-waves, positive and negative, are congruent. If, however, the third harmonic starts with maximum value coinciding with zero value of the fundamental, the resultant graph will have non-symmetrical half-waves, but successive half-waves will still be congruent. It is not difficult to see, by actual drawing, that odd-order harmonics only, superposed on a fundamental, always give congruent successive half-waves, but when any harmonic is of an even-number order a positive is followed by a negative half-wave of different shape.

Although the foregoing, a matter of mere geometry, is not

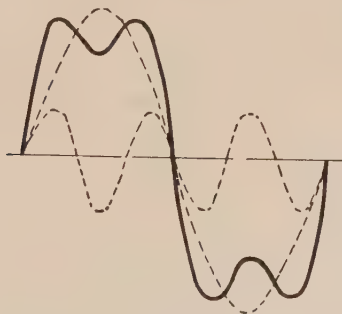


FIG. 65.

difficult to understand, Fourier's theorem is by no means self-evident, for this theorem states that any wave form can be obtained by the superposition of simple harmonic components. Thus it states that wave forms which are rectangular or triangular, can be built up more and more closely by taking more and more simple harmonic components of suitable frequencies and amplitudes. Further, the converse of Fourier's theorem is also true, that given the full particulars of a periodic function, say, in the shape of a graph showing one complete cycle of change, all particulars of the simple harmonic components, their frequencies, amplitudes, and phases, can be completely determined. This is a specialised technique of mathematics called harmonic analysis, and both

this and any kind of proof of Fourier's theorem are outside the scope of this book.

There are two important practical consequences of Fourier's theorem which can easily be established. Suppose a periodic function is equal to the sum of two harmonic components like $a \sin \omega t + b \sin (p\omega t + \theta)$. The square of this function will be the square of this sum, that is, $a^2 \sin^2 \omega t + 2ab \sin (p\omega t + \theta) \sin \omega t + b^2 \sin^2 (p\omega t + \theta)$. The average value of this square will be the sum of the average values of its components. Now, we know that the average value of the square of a simple harmonic function is half the square of its amplitude and the average of the product of two such functions of different frequencies like $2ab \sin (p\omega t + \theta) \sin \omega t$ is zero. The average square of the principal function is therefore $\frac{1}{2}(a^2 + b^2)$. A little thought will show that this result can be generalised, and that however many harmonic components there may be, the average square of the principal function will be half the sum of the squares of the amplitudes of the harmonics. The phases of the harmonics do not affect this mean square value.

Again, consider the product of two periodic functions like $a \sin \omega t$ and $b \sin (\omega t + \theta) + c \sin (p\omega t + \theta_1)$. This is $ab \sin \omega t \sin (\omega t + \theta) + ac \sin \omega t \sin (p\omega t + \theta_1)$. The average of the product is the sum of the averages of these two components. Now the product of the harmonic components of different frequencies gives a zero average, so the average product is simply the average of $ab \sin \omega t \sin (\omega t + \theta)$, which we have seen is $\frac{1}{2}ab \cos \theta$. This result can also be generalised, and the product of two periodic functions has an average value equal to the sum of such terms as $\frac{1}{2}ab \cos \theta$ where a and b are the amplitudes of harmonics of the same order, and θ is the phase difference of these equal frequency components.

Inverse Circular Functions. We have already in Chapter III come across the idea of an inverse circular function, arc $\sin x$, which means an angle the sine of which is x . As we now know the meaning of negative indices, and that $y = x^{-1}$ is a statement of what is called inverse proportionality, the reasons of the alternative symbolism, $\sin^{-1} x$ will now be clear, the -1 index indicating the inverse character of the function. With the extended idea of angle as amount of turning, or as

being proportional to time with rotation at constant angular speed, it is easy to see that the function $y = \arcsin x$ is what is called many-valued. If, for instance, $x = \frac{1}{2}$ then an angle of 30 degrees satisfies the relation, but it is also satisfied by an angle of $(180-30)$ degrees, or by a geometrical angle made by a second line symmetrical to the vertical with a 30 degree line. Further, any angle obtained by adding 1 complete turn or 360 degrees to these angles will have the same sine, viz. $\frac{1}{2}$, for complete turns merely reproduce the geometrical angle. Similar considerations apply to the arc cos and the arc tan functions. There are an unlimited number of values of the functions defined by $\arcsin x$ and $\arctan x$. The smallest value of y defined by the function, which corresponds to the angle in the geometrical sense, is called the principal value.

CHAPTER IX

DERIVED FUNCTIONS

Differential Coefficients. In Chapter VII on graphs we saw that, depending on a function, $y = f(x)$, there are two secondary functions, represented geometrically, first by the variable slope of the graph of the principal function, and, secondly, by the variable area included by the axes of co-ordinates, the graph, and the vertical ordinate through any point on it. These secondary functions were there designated respectively slope and area functions; the slope function was defined as the limiting value of the average slope between two neighbouring points of the graph; and the principal function was found to stand in the same relation to the area function, as the slope function does to the principal function. The area function can thus, in a sense, be considered to be an inverse slope function. Again, in Chapter VIII the idea of a secondary function was encountered in the consideration of the continually changing velocity of a point executing a periodic simple harmonic motion, and it was there seen that the slope function of a graph showing the variation in time of a moving point is proportional to the function expressing instantaneous speed in terms of time. We have now to consider these ideas of secondary functions in a general and abstract way, and apart from the geometrical idea of the slope of a graph or the kinematic idea of a physical speed or velocity.

The mathematical statement $y = f(x)$ means that assigning any value to x gives a certain value or values to y . If x be considered for the moment to represent some definite number, then if this increases by a small amount from x to $(x + \delta x)$, y will change from $f(x)$ to $f(x + \delta x)$, and the change in y , or δy , will, if we can evaluate both $f(x)$ and $f(x + \delta x)$, be $f(x + \delta x) - f(x)$.

The ratio $\frac{\delta y}{\delta x}$ gives a kind of average rate at which y changes with x in the interval δx . If this ratio tends to a definite limiting value as δx diminishes, that is, if the ratio can be made to approximate to a limiting value as closely as may be

desired, by diminishing δx , then this limiting value of the ratio will, as it is calculated from $f(x) - f(x + \delta x)$ and δx , be a new function of x . It is called the first derived function of $f(x)$ and denoted by $f'(x)$. It is also called the differential coefficient of y with respect to x , denoted by the symbol $\frac{dy}{dx}$ (read dy by dx) and the determination of the nature of this secondary or derived function is a mathematical process called differentiation, which is the subject matter of that branch of mathematics called the differential calculus.

Let us return to our algebraic illustration on p. 121. If $y = x^2$, this expresses a definite functional relationship between y and the independent variable x . The values of y corresponding to x and $x + \delta x$ are respectively x^2 and $x^2 + 2x\delta x + (\delta x)^2$, so that the change of y , δy , corresponding to the change of x , δx , is $2x\delta x + (\delta x)^2$, and the ratio $\frac{\delta y}{\delta x} = 2x + \delta x$.

This is true so long as δy and δx have any values however small, but it has no meaning when these values are each zero. Suppose that $x = 1$ and the δx is one-millionth, or as we can

express it, 10^{-6} . We see that δy is $2 \times 1 \times 10^{-6} + 10^{-12}$, and $\frac{\delta y}{\delta x}$ has the definite value of $2 + 10^{-6}$. We can make δx smaller still, say, 10^{-10} , in this case $\frac{\delta y}{\delta x}$ has the value, $2 + 10^{-10}$.

Thus the smaller we make δx the closer does the value of $\frac{\delta y}{\delta x}$, or $2x + \delta x$ approximate to the value $2x$, in this case 2. We can conclude, therefore, that the limiting value of $\frac{\delta y}{\delta x}$ is equal to $2x$, and this is expressed mathematically by the statement that if $y = x^2$, $\frac{dy}{dx} = 2x$. It is easy to see that if $y = ax^2$, δy will be increased a -fold, and that $\frac{dy}{dx}$ will be $2ax$.

Before we go any further we must pause to consider this new symbol $\frac{dy}{dx}$. We have defined it as the limiting value of

a ratio, and it looks like a ratio that of (dy) to (dx) . But we cannot say that (dy) or (dx) stand for any quantities, certainly not for zero values of δy and δx , for the fraction $\frac{0}{0}$ has no value.

Sometimes $\frac{dy}{dx}$ is considered as a symbol attaching to the func-

tion y , and on this view it can be written $\frac{d}{dx}(y)$. Thus if the

function y is $\arccos x$, $\frac{dy}{dx}$ might be written $\frac{d}{dx}(\arccos x)$,

the $\frac{d}{dx}$ symbol denoting something derived from y . But,

as we shall see later, the two terms of the apparent ratio, dy and dx , may be used standing alone as if they represented actual numbers or magnitudes, which they certainly do not. There

is no doubt that this $\frac{dy}{dx}$ symbol is unfortunate and equivocal.

Primarily it looks like $d \times y$ divided by $d \times x$, and even if dy and dx are reckoned single symbols it looks like the ratio of

two definite quantities. The $\frac{dy}{dx}$ notation is firmly established

in mathematics. As it stands it must be considered a complete symbol for the first derived function of x corresponding

to $y=f(x)$. In the formula $\frac{d}{dx}f(x)$, $\frac{d}{dx}$ denotes a particular

calculation called differentiation carried out on $f(x)$. We

shall refer later to the meaning of the symbols dy and dx when they are found separately. Whoever desires to under-

stand the language of mathematics must be familiar with the

$\frac{dy}{dx}$ notation, and with the several and not altogether consistent

meanings that are, by convention, given to it.

The Technique of Differentiation. It appears from our definition of a differential coefficient that if $y=f(x)$ and we can interpret $f(x+\delta x)$ algebraically, then we can calculate the change δy of y corresponding to the change δx of x , and the value to which the ratio of δy to δx approximates the more closely, as δx is diminished, is the required value of the derived function $\frac{dy}{dx}$. This method of differentiation is the

basic one, and it is used for the functions of a fundamental character. The actual process of differentiation can in practice usually be simplified by utilising certain important properties of differential coefficients whereby the differentiation of a function can be made to depend upon the known differential coefficients of simpler component functions.

As an example of the basic and direct process of differentiation we may consider the calculation of $\frac{dy}{dx}$ for $y = x^n$, when n is a positive integer or whole number. It follows from the argument on p. 23 that the evaluation of the expression $(x + \delta x)^n$ by algebra gives a result containing decreasing powers of x , the first term is x^n or y , the second is $nx^{n-1} \times \delta x$, the third contains the quantity $(\delta x)^2$, and so on. It follows therefore that when $y = x^n$, $\delta y = (x + \delta x)^n - x^n = nx^{n-1}\delta x$ plus terms containing the square and higher powers of δx . $\frac{\delta y}{\delta x}$ is therefore equal to nx^{n-1} plus terms containing δx and its powers. The smaller δx is the smaller still will be its powers, so that, by the argument previously employed, it follows that when $y = x^n$, n being a positive integer, $\frac{dy}{dx} = nx^{n-1}$.

We shall deal with the simplest properties of differential coefficients used in the technique of differentiation; for a full account of these properties any standard textbook on the calculus may be consulted. Suppose y is the product of two functions of x , say $u = f(x)$ and $v = \psi(x)$. A change δx in x will turn u into $u + \delta u$ and v into $v + \delta v$ so that the $y + \delta y$ will be $(u + \delta u)(v + \delta v)$. Working out this product we find that $y + \delta y = uv + u\delta v + v\delta u + \delta u\delta v$, and as $uv = y$ we find that $\frac{\delta y}{\delta x} = u\frac{\delta v}{\delta x} + v\frac{\delta u}{\delta x} + \delta u\frac{\delta v}{\delta x}$. By diminishing δx we can make the change of u or δu in the last term as small as we desire so that the limiting value of $\frac{\delta y}{\delta x}$ is equal to the sum of the limiting

values of the terms in its equivalent, and $\frac{dy}{dx} = u\frac{dv}{dx} + v\frac{du}{dx}$.

Let us illustrate this. If $y = x^6$ we know, by the rule found

by the direct method, that $\frac{dy}{dx} = 6x^5$. But y is also equal to $x^4 \times x^2$, so that if $u = x^4$ and $v = x^2$, $\frac{du}{dx}$ is $4x^3$ and $\frac{dv}{dx} = 2x$, and using the product of functions rule for differentiation $\frac{dy}{dx} = x^4 \times 2x + x^2 \times 4x^3 = 6x^5$, as we found by the direct rule.

Suppose that y is a function of a function of x , or $y = \psi(u)$ where $u = f(x)$. If x changes by δx , this change produces a change δu in u , and this change δu changes y by δy , then $\frac{\delta y}{\delta x} = \frac{\delta y}{\delta u} \times \frac{\delta u}{\delta x}$, for the two δu 's on the right-hand side of this equation refer to the same quantity. This equation is true for the limiting values of the fractions or for the corresponding differential coefficients so that $\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx}$.

This can be illustrated very simply. If $y = x^6$, y is also equal to $(x^3)^2$, that is, to the square of the function x^3 . Representing x^3 by u , we have, as $y = u^2$, $\frac{dy}{dx} = 2u \times 3x^2$, which, as $u = x^3$ is equal to $6x^5$, as was found by two other methods.

If y is a function of x , then x is a function of y , and $1 \div \frac{\delta y}{\delta x}$ or $\frac{\delta x}{\delta y}$ gives the ratio of the change in x corresponding to a change in y . Thus the limiting value of $\frac{\delta x}{\delta y}$ or $1 \div \frac{dy}{dx}$ is the differential coefficient of x considered as a function of y .

Differentiation of x^n . We have differentiated the function x^n , where n is a positive whole number and found that $\frac{dy}{dx} = nx^{n-1}$. We can use the results established in the preceding section to show that this formula is correct without this restriction regarding n .

Suppose $x = y^m$ where m is a positive integer, this of course means that $y = x^{\frac{1}{m}}$ or that y is the m th root of x . We have

$\frac{dx}{dy} = my^{m-1}$ or $m \frac{y^m}{y}$, which is equal to $m \frac{x}{y}$. Thus $\frac{dy}{dx} = \frac{y}{mx}$ or $\frac{1}{m} \times \frac{x^{\frac{1}{m}}}{x}$ or to $\frac{1}{m} x^{\frac{1}{m}-1}$, so that if n , instead of being a positive integer is a positive fraction with 1 as a numerator the $\frac{dy}{dx}$ of x^n is still nx^{n-1} .

If $y = x^{\frac{m}{l}}$, m and l being positive integers, we can by treating $x^{\frac{m}{l}}$ as being equal to $(x^{\frac{1}{l}})^m$, and using the rule for the differentiation of a function show that the $\frac{dy}{dx}$ of x^n is still nx^{n-1} when n is any positive fraction.

Suppose $y = x^{-n}$ or $\frac{1}{x^n}$, then $yx^n = 1$; yx^n is constant and its differential coefficient for its change when x changes must be zero. Let us differentiate yx^n as a product of two functions. The result is $y \times nx^{n-1} + x^n \frac{dy}{dx}$ or as $y = x^{-n}$, $x^{-n} \times nx^{n-1} + x^n \frac{dy}{dx} = \frac{n}{x} + x^n \times \frac{dy}{dx}$, and this sum is equal to zero. It follows,

therefore, that $\frac{dy}{dx} = -n \times \frac{1}{x^{n+1}}$ or to $-nx^{-n-1}$, so that the rule is correct when n is a negative number integral or fractional. Although we have not proved the rule when n is an irrational number, and this proof is beyond the scope of this book, we may fairly conclude that if n is any real number, and $y = x^n$ then $\frac{dy}{dx} = nx^{n-1}$. The foregoing reasoning, although a little involved, is well worth following carefully as it shows how, by using the basic properties of differential coefficients, we can generalise a formula like $\frac{dy}{dx} = nx^{n-1}$ which we have

established by algebra to be true when n is a positive integer, and remove the restriction of the nature of the exponent n .

Let us now consider the results of differentiating some powers of x by the general formula as set out in the following scheme :

$$y = x^3 \quad \frac{dy}{dx} = 3 \times x^{3-1} = 3x^2.$$

$$y = x^2 \quad \frac{dy}{dx} = 2 \times x^{2-1} = 2x.$$

$$y = x^1 = x \quad \frac{dy}{dx} = 1 \times x^{1-1} = 1 \times x^0 = 1.$$

$$y = x^0 = 1 \quad \frac{dy}{dx} = 0 \times x^{0-1} = 0.$$

$$y = x^{-1} \quad \frac{dy}{dx} = -1 \times x^{-1-1} = -x^{-2}.$$

$$y = x^{-2} \quad \frac{dy}{dx} = -2 \times x^{-2-1} = -2x^{-3}.$$

In this list the powers of x in the principal functions represented by x^n diminish regularly from 3 to -2 , the powers of x in the differential coefficients also diminish from 2 to -3 . We notice, however, a remarkable thing: every power of x is represented in the principal functions, but there is no -1 power in the differential coefficients, in short, no power of x when differentiated gives x^{-1} multiplied by a constant number for an answer. It looks as if x^0 is the power that ought to give a differential coefficient containing x^{-1} , but according to the formula, the multiplier of x^{-1} must be zero. We know, of course, that if $y = x^0 = 1$, y is constant and, as it cannot change however x is varied, δy , and hence $\frac{dy}{dx}$ must be zero.

Logarithmic and Exponential Functions. Let us assume what seems reasonable, that there is some function y of x which has the property that $\frac{dy}{dx}$ is equal to x^{-1} or $\frac{1}{x}$. We shall call this function $\log x$. At present $\log x$ is merely a name of a function of which we know nothing excepting the value of its differential coefficient. If $y = \log x$ and $\frac{dy}{dx} = \frac{1}{x}$, then we know that $\frac{dx}{dy} = x$. x is an inverse function of $y = \log x$, or, using x and y in their usual significance as standing for independent and dependent variables respectively, if there is

a function y of x such that $\frac{dy}{dx} = \frac{1}{x}$, there is also a function y of x such that $\frac{dy}{dx} = y$. We call this latter function $\exp x$.

When a function is defined by an equation like $\frac{dy}{dx} = \frac{1}{x}$ for $\log x$, or like $\frac{dy}{dx} = y$ for $\exp x$, such an equation is called a differential equation. These two equations tell us very little explicitly about the functions $\log x$ and $\exp x$. We shall, however, see later how we can, by the differential calculus, discover a little more about the $\exp x$ function, but we shall have to defer a complete elucidation of this, and also of the $\log x$ function till the next chapter. At present they are functions, the existence of which we assume and which, if they exist, satisfy certain differential equations.

Differentiation of Circular Functions. We have already in Chapter VIII worked out by the direct method the differential coefficients of the functions $\sin \theta$ and $\cos \theta$, for when $x = a \sin \omega t$ the limiting value of the velocity corresponding to a displacement x is $a\omega \cos \omega t$, and this velocity is, according to our definition of a differential coefficient, the result of differentiating the function $a \sin \omega t$ with respect to t . If $\omega = 1$ then t is the measure of the angle ωt which may be called θ , so that if $x = a \sin \theta$ $\frac{dx}{d\theta} = a \cos \theta$. Similarly, when investigating the rate of change of the velocity $a\omega \cos \omega t$, we arrived at a result equivalent to $\frac{dx}{d\theta} = -a \sin \theta$ when $x = a \cos \theta$.

Let us consider the function $x = a \tan \theta$. We know from Chapter III that $\tan \theta = \frac{\sin \theta}{\cos \theta}$ so that this equation is equivalent to $x = a \frac{\sin \theta}{\cos \theta}$, or $x \cos \theta = a \sin \theta$, and the differential coefficients of the two sides of this equation must be equal. We can find the differential coefficient of the left-hand side by the rule for the product of functions. It is $x \times -\sin \theta + \cos \theta \times \frac{dx}{d\theta}$, and this equals $a \cos \theta$. Dividing this last

equation throughout by $\cos \theta$ we find that $-x \times \frac{\sin \theta}{\cos \theta} + \frac{dy}{dx} = a$, or as $x = a \tan \theta$, and $\frac{\sin \theta}{\cos \theta} = \tan \theta$, $-a \tan^2 \theta + \frac{dy}{dx} = a$. We thus conclude that when $x = a \tan \theta$, $\frac{dx}{d\theta} = a(1 + \tan^2 \theta)$.

We can readily obtain the differential coefficients of the inverse circular functions. For if $\theta = \arcsin \frac{x}{a}$, this means that $x = a \sin \theta$ so that $\frac{dx}{d\theta} = a \cos \theta$ and $\frac{d\theta}{dx} = \frac{1}{a \cos \theta}$. But, as $\cos^2 \theta + \sin^2 \theta = 1$, $\cos^2 \theta = 1 - \sin^2 \theta = 1 - \frac{x^2}{a^2} = \frac{a^2 - x^2}{a^2}$, so that $a \cos \theta = \sqrt{(a^2 - x^2)}$. Thus $\frac{d\theta}{dx} = \frac{1}{\sqrt{(a^2 - x^2)}}$ when $\theta = \arcsin \frac{x}{a}$. We arrive at a similar result if we calculate the differential coefficient of $\arccos \frac{x}{a}$. Now consider $\theta = \arctan \frac{x}{a}$ or $x = a \tan \theta$. As $\frac{dx}{d\theta} = a(1 + \tan^2 \theta)$ or $a \left(1 + \frac{x^2}{a^2}\right) = \frac{a^2 + x^2}{a}$ we find that $\frac{d\theta}{dx} = \frac{a}{a^2 + x^2}$. We ought to note that whether θ is defined as $\arcsin \frac{x}{a}$, $\arccos \frac{x}{a}$, or $\arctan \frac{x}{a}$, the ratio $\frac{\delta \theta}{\delta x}$ has the dimension of a ratio or pure number divided by a length, so that the differential coefficient of θ in each case ought to have the dimensions of 1 divided by length, as it has.

The reader ought also to notice that whereas the differentiation of the circular functions gives other circular functions, the differentiation of inverse circular functions gives algebraic functions, which are of a much simpler character. He ought to think carefully why this is so, and realise clearly that it depends ultimately upon Pythagoras's theorem of geometry.

Successive Differentiation. The differential coefficient $\frac{dy}{dx}$ of $y = f(x)$ is a new function of x which we have already called the first derived function and indicated by the symbol $f'(x)$

This new function can generally itself be differentiated with respect to x . Thus, if $y = x^3 = \frac{dy}{dx} = 3x^2$, and the differential coefficient of the derived function $3x^2$ is $6x$, the result of a second differentiation like this is called the second derived function denoted by $f''(x)$, or, more usually, the second differential coefficient, denoted by the symbol $\frac{d^2y}{dx^2}$.

With all but algebraic integral functions this successive differentiation can be carried on indefinitely. The result of the n th differentiation is denoted by $\frac{d^n y}{dx^n}$. With an integral function composed of terms like ax^n where n is a positive integer the process is limited, for as each differentiation diminishes the index of x by 1, a differential coefficient equal to x multiplied by some constant will ultimately be obtained. The next differential coefficient will be this constant, and after this all others will be zero. If, however, the index n of ax^n is a negative or fractional number the successive differentiation can proceed indefinitely; for example, if $y = x^{\frac{1}{2}}$ the indices of the powers of x with successive differentiations will be $-\frac{1}{2}$, $-\frac{3}{2}$, $-\frac{5}{2}$, and so on. The reader who pursues his studies of mathematics will find that this is connected with the fact that the expression $(x+a)^n$ can be expressed as an integral function with a limited number of terms if n is a positive whole number, but if n is fractional or negative $(x+a)^n$ can only be changed into an infinite series.

The successive differentiation of $\exp x$ is peculiar, for as $\frac{dy}{dx} = y$ each successive differentiation merely reproduces the original function.

Successive differentiation of $y = a \sin \theta$ gives an interesting result, for as $\frac{dy}{d\theta} = a \cos \theta$ and $\frac{d^2y}{d\theta^2} = -a \sin \theta$, two differentiations reproduce the function with a change of sign. Four differentiations will correct this change of sign and repeat the original function.

We have seen on p. 158 that the simple harmonic motion $x = a \sin \omega t$ can be represented by a directed line or vector of

length a ; $\frac{dx}{dt}$ being $a\omega \cos \omega t$ can be represented by a second vector 90 degrees in advance of the first, and of length $a\omega$. Differentiating $a \sin \omega t$ is therefore equivalent to multiplying the complex number representing the corresponding vector by $j\omega$. The second differentiation of $a \sin \omega t$ gives $a\omega^2 \sin \omega t$, vector representing this second derived function is ω^2 times the primary vector in length, and reverse to it in direction. The complex number representing the second derived function vector is therefore obtained by multiplication of the primary complex number by $-\omega^2$ or by $j\omega \times j\omega$. It follows, therefore, that the operation of differentiation represented by $\frac{d}{dt}$ applied to simple harmonic functions like $a \sin \omega t$, is equivalent to multiplication of a number by $j\omega$, and when $\omega = 1$, $\frac{d}{dt}$ is equivalent in this sense to j . This applies, of course, only to circular or simple harmonic periodic functions, and the exploitation of this property of the $\frac{d}{dt}$ symbol of differentiation, whereby it is in a way equivalent to an algebraic operation, is the foundation of a very powerful technique of higher mathematics called the calculus of operations.

Geometrically, as $\frac{dy}{dx}$ is the slope of the graph representing the functional relationship of y to x , so $\frac{d^2y}{dx^2}$ is the limiting value of the linear rate at which the slope changes as the x ordinate of the curve increases. We have seen on p. 125 that where a graph shows a maximum value of y the slope changes from upward or positive to downward or negative. As the slope is decreasing where the graph passes through a maximum value, the value of $\frac{d^2y}{dx^2}$ corresponding will be negative. Similarly, the sign of $\frac{d^2y}{dx^2}$ corresponding to a minimum value will be positive. These statements are generally true, although it is possible to have maximum or minimum values with the

second derivative $\frac{d^2y}{dx^2}$ equal to zero, like the first derivative or slope. The discussion of these anomalous cases is outside the scope of this book.

Turning Values, Maxima and Minima. We have already come across the idea of a turning value of a function, a point at which the slope of its graph changes from up to down or vice versa, and we have seen that a turning value gives either a local maximum or minimum value to the function, and that the geometrical criterion for a turning value is that the slope of the graph is zero. This criterion is the same as $\frac{dy}{dx} = 0$.

Thus, if we can calculate the $\frac{dy}{dx}$ of a function $y = f(x)$ then, by equating $\frac{dy}{dx}$ to zero, we obtain an equation with x as an un-

known, the solutions of which, if they can be obtained, will give the x values corresponding to the maximum or minimum values. We saw, moreover, in the preceding section that we can, having found a turning value in this way, generally determine whether it is a maximum or minimum, by calculating the second derived function $\frac{d^2y}{dx^2}$ and finding whether it is positive

or negative for the critical value or values of x that give the turning values. The determination of the maximum and minimum values of functions is a very important application of the differential calculus, which is fully treated in the ordinary textbooks. We may note that the mechanism of this technique does not always give us quite the answer we want. Thus consider the function $ax^5 + bx^3 + cx$. We can easily find the differential coefficient; it is $6ax^4 + 3bx^2 + c$, and, although we can say that the values of x that give turning values to the function are those that satisfy the equation $\frac{dy}{dx} = 0$, that is, $6ax^4 + 3bx^2 + c = 0$, we cannot say what these values are, because we cannot solve this equation in a formal way, and, hence, we cannot obtain the turning values in terms of the constants a , b , and c in the function.

A simple example of the discovery of a turning value is

given by the function $y = ax + \frac{b}{x}$, $\frac{dy}{dx} = a - \frac{b}{x^2}$, and if $\frac{dy}{dx} = 0$ $x = \sqrt{\frac{b}{a}}$, and substituting this in the function we find a turning value of y to be $2\sqrt{(ab)}$. The second derived function is obtained by differentiating $a - \frac{b}{x^2}$, this gives $-\left(\frac{-2b}{x^3}\right) = \frac{2b}{x^3}$, and when $x = \sqrt{\frac{b}{a}}$, and is positive, the value of $\frac{2b}{x^3}$ is also positive, so that the turning value is a minimum.

Again, consider the function $x = \sin \omega t + \cos \omega t$, $\frac{dy}{dt} = \omega \cos \omega t - \omega \sin \omega t$, and this is zero when $\cos \omega t = \sin \omega t$ or when $\tan \omega t = 1$. The angle that has 1 for its tangent is $\frac{\pi}{4}$ or 45 degrees, the sine and cosine of which are both $\frac{1}{\sqrt{2}}$. Thus the turning value of the function is $\frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} = \sqrt{2}$. We see that $\frac{d^2y}{dt^2} = -\sin \omega t - \cos \omega t$, and when the sine and cosine are each $\frac{1}{\sqrt{2}}$ positive, the second derivative is negative, hence the turning value is a maximum.

As a matter of fact we need not have gone to all this trouble. We knew that the two functions $\sin \omega t$ and $\cos \omega t$ are simple harmonic, of equal amplitude, whereof $\cos \omega t$ is 90 degrees or $\frac{\pi}{2}$ out of phase with $\sin \omega t$. These functions can, as was explained on p. 151, be represented by vectors at 90 degrees, each equal to 1, and the sum of the functions is a simple harmonic function represented by the resultant of these two vectors, that is, by a vector of length $\sqrt{2}$ bisecting the angle between them. This resultant is, therefore, $\sqrt{2} \sin \left(\omega t + \frac{\pi}{4} \right)$, and its maximum value is evidently $\sqrt{2}$. when the sine is unity, and $\omega t + \frac{\pi}{4} = \frac{\pi}{2}$.

Similarly, as $a \sin \omega t + b \cos \omega t$ can be represented graphically as shown in Fig. 66 by two vectors a and b at right angles, this function is equivalent to a simple harmonic function having an amplitude corresponding to the resultant vector, i.e. $\sqrt{a^2 + b^2}$, and leading $\sin \omega t$ by an angle $\phi = \arctan \frac{b}{a}$. Thus, as $a \sin \omega t + b \cos \omega t = \sqrt{a^2 + b^2} \sin (\omega t + \phi)$ its maximum value must be $\sqrt{a^2 + b^2}$ when $\omega t + \phi = \frac{\pi}{2}$.

The two preceding paragraphs illustrate the fact that, although the differential calculus provides a kind of mechanical technique for the discovery of the turning values of a function,

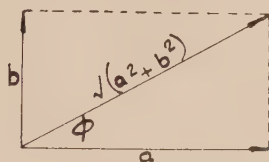


FIG. 66.

this method is not always the quickest, the best, or the most illuminating. As a matter of fact, while the calculus give the condition for the turning value of a function, that its differential coefficient is zero, the trigonometrical method of the preceding paragraph gives, not only this condition, but also the maximum value of the function. Similarly, by converting the function referred to at the top of page 196 into the sum of two squares by the artifice described on page 150, both the condition and also the minimum value are at once revealed.

Expansion of Functions. We shall illustrate by one or two very simple examples a more subtle and far-reaching use of the differential calculus, that of the expansion of functions into infinite series. We have already pointed out when dealing with the basic idea of function that, to evaluate any but algebraic functions, such an expansion into a converging series, is essential for numerical evaluation. Thus, the func-

tion $\sin \theta$ is geometrically defined as the relation between a particular side of a right-angled triangle of unit hypotenuse and an arc of a circle associated with this triangle in a particular way. The geometrical definition give no other way of determining numerical values of $\sin \theta$ than by actual measurement of the lines of a geometrical figure. Now, we know from geometrical considerations, several properties of the $\sin \theta$ function: (a) that it is continuous, (b) that $\sin \theta = 0$ when $\theta = 0$, (c) that when θ is very small $\sin \theta$ is approximately equal to the circular measure of θ , and (d) that $\sin \theta$ is an odd function, that is to say, $\sin(-\theta) = -\sin \theta$. Lastly, the differential calculus tells us that if $\sin \theta$ is differentiated twice with respect to θ the result is $-\sin \theta$. This information is sufficient to enable us to obtain an infinite series representing $\sin \theta$, and from which assigning a value to θ , $\sin \theta$ can be calculated arithmetically and independently of any geometrical construction.

Let us assume that $\sin \theta$ is identically equal to an infinite series of the form $a + b\theta + c\theta^2 + d\theta^3 \dots$, that is, that the series is true for all values of θ . Then, as $\sin \theta = 0$ when $\theta = 0$ we know at once that a in the series is 0; if it had any other value $\sin \theta$ would be a when θ is 0. The series must thus be like $a\theta + b\theta^2 + c\theta^3 \dots$. Again, as $\sin \theta$ is an odd function we know it can contain only odd powers of θ , for it is only these powers which will change sign with a change of the sign of θ . The series must therefore be of the form $a\theta + b\theta^3 + c\theta^5 \dots$. Further, if θ is expressed arithmetically in circular or radian measure, then when θ is very small so that terms in the series containing θ^3 and higher powers can to a first approximation be neglected, $\sin \theta = a\theta$, and a must be equal to 1, to give the geometrical approximation $\sin \theta = \theta$. We therefore find that if there is an infinite series for $\sin \theta$ it must have the form $\sin \theta = \theta + a\theta^3 + b\theta^5 \dots$, and we have to discover what numbers the letters a , b , etc., stand for. It is at this point that the differential calculus comes into the investigation, for the second differential coefficient of $\sin \theta$ must be equal to the second differential coefficient of the series to which it is equivalent, and this series can easily be differentiated term by term, by using the x^n rule. We have therefore

$$y = \sin \theta = \theta + a\theta^3 + b\theta^5 + c\theta^7 + d\theta^9 \dots$$

$$\frac{dy}{d\theta} = \cos \theta = 1 + 3a\theta^2 + 5b\theta^4 + 7c\theta^6 + 9d\theta^8 \dots$$

$$\frac{d^2y}{d\theta^2} = -\sin \theta = 3 \times 2a\theta + 5 \times 4b\theta^3 + 7 \times 6c\theta^5 + 9 \times 8d\theta^7 \dots$$

$$\text{and } +\sin \theta = -3 \times 2a\theta - 5 \times 4b\theta^3 - 7 \times 6c\theta^5 - 9 \times 8d\theta^7 \dots$$

We have here two series, the first and fourth, both of which are equal to $\sin \theta$, for all values of θ . By the principle of undetermined coefficients explained on p. 144, it follows that the coefficients of identical powers of θ in the two series must be equal. For instance, the coefficient of θ in the first series is 1, in the fourth series it is $-3 \times 2a$, so $-3 \times 2a = 1$ and $a = -\frac{1}{3 \times 2}$ or, using the factorial notation $a = -\frac{1}{3!}$. Similarly,

dealing with θ^3 terms, we see that $a = -5 \times 4b$, or $-\frac{1}{3 \times 2}$

$= -5 \times 4b$ so that $b = \frac{1}{2 \times 8 \times 4 \times 5}$ or $\frac{1}{5!}$. A little thought

shows that, by using this process, each letter coefficient in the first series comes out to be numerically equal to 1 divided by the factorial of the power of θ with which it is associated, and that the numerical coefficients are alternately positive and negative. Thus we find that

$$\sin \theta = \theta - \frac{\theta^3}{3!} + \frac{\theta^5}{5!} - \frac{\theta^7}{7!} \dots$$

We have seen on p. 98 that if we arrive at our infinite series that purports to give a means for the numerical evaluation of a function we must satisfy ourselves that this series is convergent. Let us apply the test of p. 103; any term, say the n th, is obtained by multiplying the preceding term by $\frac{\theta^2}{n(n-1)}$,

and this factor can be made as small as we please, however large θ may be, by sufficiently increasing n . The series would, therefore, be convergent for all values of θ if the signs of all the terms were positive. It is, a fortiori, convergent with the signs alternating.

We can, by a similar process, work out a converging infinite series for $\cos \theta$. We know that when $\theta = 0$, $\cos \theta = 1$, that $\cos \theta$

is an even function, so that a series representing it can contain only even powers of θ , and that two differentiations of $\cos \theta$ give $-\cos \theta$. We can therefore write a symbolic series for $y = \cos \theta$ as

$$y = \cos \theta = 1 + a\theta^2 + b\theta^4 + c\theta^6 \dots$$

$$\frac{dy}{d\theta} = -\sin \theta = 2a\theta + 4b\theta^3 + 6c\theta^5 \dots$$

$$\frac{d^2y}{d\theta^2} = -\cos \theta = 2a + 4 \times 3b\theta^2 + 6 \times 5c\theta^4 \dots$$

$$\text{and } \cos \theta = -2a - 4 \times 3b\theta^2 - 6 \times 5c\theta^4.$$

Equating coefficients of identical powers of θ in the first and fourth series we find $1 = -2a$ so that $a = -\frac{1}{2}$ or $-\frac{1}{2!}$,

$a = -4 \times 3b = -\frac{1}{2}$ so that $b = +\frac{1}{2 \times 3 \times 4} = \frac{1}{4!}$, and so on. Thus

$$\cos \theta = 1 - \frac{\theta^2}{2!} + \frac{\theta^4}{4!} - \frac{\theta^6}{6!} \dots$$

a series which can easily be seen to be convergent for all values of θ .

We must realise clearly what these series mean. If a value of an angle θ is assigned numerically in circular measure, then substituting this value of θ in the series we ought to be able to obtain a numerical value for either $\sin \theta$ or $\cos \theta$. The numerical value of θ must be in radians, for the property $\sin \theta = \theta$ as θ approaches 0 as a limit only holds for this condition. Let us try how the series for $\sin \theta$ works out for an angle of 1 radian by putting $\theta = 1$ in the series. We may set out the sum in the following way :

	+	-	
$3 \times 2 = 6$)1.00000	1.00000		
$4 \times 5 = 20$)0.16666		0.16666	
6)0.00833	0.00833		
7)0.00139			
8)0.00020		0.00020	
9)0.00002			
0.00000			
	1.00833	0.16686	
	-0.16686		
	0.84147		

The series therefore gives a value of 0.84147 for the sine of an angle of 1 radian. This angle is 57.3 degrees very nearly, and the reader ought to be sufficiently interested to look up the sine of 57.3 degrees in a book of reference and see if the value we have found is approximately correct.

There is another point about the series that a thinking reader will question. $\sin \theta$ increases from 0 to 1 as θ increases from 0 to $\frac{\pi}{2}$, and then decreases from 1 to 0 as θ further increases from $\frac{\pi}{2}$ to π . If therefore $\theta = \pi = 3.14159$ be substituted in the series formula the arithmetical result of calculating $\sin \theta$ ought to be zero. Will this be so? It can be shown by higher algebra that the equation $0 = \theta = \frac{\theta^3}{3!} + \frac{\theta^5}{5!} \dots$ is satisfied by an infinite number of positive roots all of which are multiples of a number just over 3, and we can be quite sure that, if a numerical value of π were substituted in the series and it were evaluated numerically, the result would approximate the closer and closer to zero as our numerical value for the number π increased in approximation.

As a final example of the expansion of a function into an infinite series we may consider $y = \exp x$, which we have defined by the differential equation $\frac{dy}{dx} = y$. If we assume, and this we must note is an assumption, that $\exp 0 = 1$ we can write a symbolic series thus

$$y = \exp x = 1 + ax + bx^2 + cx^3 + dx^4 \dots$$

$$\text{and } \frac{dy}{dx} = \exp x = a + 2bx + 3cx^2 + 4dx^3 \dots$$

so that equating coefficients of like powers of x in the two series for $\exp x$ we find $a = 1$, $2b = a = 1$ and $b = \frac{1}{2}$, $b = 3c = \frac{1}{2}$ and $c = \frac{1}{2 \times 3} = \frac{1}{3!}$, and so on. Thus

$$\exp x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} \dots$$

A series identical with that given on p. 98 as equal to a function e^x . We gave the e^x series merely as an illustration;

we have deduced the $\exp x$ series from the definitions (a) that if $y = \exp x$ $\frac{dy}{dx} = y$, and (b) $\exp 0 = 1$. We are not justified in saying that the series is equal to a number raised to a variable power. All we know, at present, is that the series we have obtained has the properties of the function $\exp x$ that we have postulated.

Integrals. As every continuous function has a differential coefficient, so it has what in geometrical language has been called an area function. If $y = f(x)$ then the algebraic representation of the area function of y is called the integral of y with respect to x , and it is denoted by the symbol $\int y dx$. This symbol requires some explanation. Consider Fig. 67 in which BA_2A_1 is the graph of a function $y = f(x)$. The value of the area function corresponding to the point A_1 will be the geometrical area denoted by OBA_1x_1 . If this area is divided, as shown, into a number of vertical strips, all of equal width, δx , then the area will be approximately equal to the sum of component areas, each of which will be $\delta x \times$ average height of the strip, or $\delta x \times y$ ordinate at the middle of the strip. A sum of this kind, that is, a sum of a number of quantities all of which can be expressed in the same way, is algebraically expressed as $\sum y \times \delta x$, Σ , the Greek letter "sigma," the equivalent of S, denoting the sum. Now it is easy to see that, although a formula of this kind represents the area OBA_1x_1 only approximately, for it assumes that the upper termination of each strip is straight, yet, the smaller the width δx is made, and the larger the number of strips are used, the closer will the algebraic sum approximate to the actual area. The value to which $\sum y \times \delta x$ can be made to approximate as closely as we please by increasing the number of strips and reducing their width is the limiting value of $\sum y \times \delta x$ as δx approaches zero, and it is indicated by the integral $\int y dx$. As we showed on p. 122 the differential coefficient of this integral is y .

Here we must note an important point. Suppose, in Fig. 67, we had reckoned x ordinates, not from the point O, but from another point O_1 to the left of it. The area function corresponding to the point A_1 would then have been $O_1B_1A_1x_1$,

a larger quantity, but the differential coefficient of this function would still be the same y corresponding to the vertical ordinate at the point x_1 . Thus the area function corresponding to a variable point on the graph is not fixed ; it depends upon the datum from which we reckon x ordinates. Otherwise, if

A is an area function or an integral such that $\frac{dA}{dx} = y$, then if we increase or decrease A by any constant quantity C the differential coefficient of the new function A + C will still be y, for the differential coefficient of a constant is zero. Thus

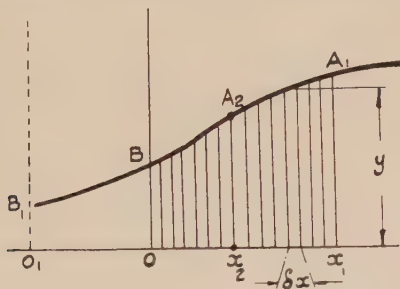


FIG. 67.

$\int y dx + C$ is the integral of y as well as $\int y dx$, and owing to this ambiguity $\int y dx$ is called an indefinite integral.

Suppose now that we consider the area under a definite section of the curve of Fig. 67, like that defined by two points A_2 and A_1 corresponding to ordinates x_2 and x_1 . Now, from whatever datum we measure these x values the area from the left up to x_2A_2 will be obtained by using the value of x_2 in the function $\int ydx + C$, but we shall not be able to evaluate this numerically because the constant C is unknown. The area up to the ordinate x_1A_1 will be obtained by substituting the value x_1 in the integral, and this will consist of a determinate number plus the same unknown constant. The area we require, $x_2A_2A_1x_1$, is the difference of these two areas; it is $\int ydx + C$ evaluated, as far as we can, for x_1 less $\int ydx + C$ evaluated for x_2 . The unknown constant disappears in the subtraction and the difference is denoted by the symbol $\int_{x_2}^{x_1} ydx$, which means

that the function $\int y dx$ is to be evaluated, first for x_1 , then for x_2 , and the second value is to be subtracted from the first.

$\int_{x_2}^{x_1} y dx$ is called a definite integral, because, giving numerical values to x_1 , and x_2 to the constant coefficients in the function $\int y dx$, its numerical value can be worked out.

Integration. The mathematical process of obtaining an integral of a given function is called integration. There is no straightforward set of rules for this process comparable in simplicity to those applicable to differentiation. The technique of integration consists inherently in the manipulation of functions so that they are turned into forms or separated into components that are known to be differential coefficients. This is not always possible, and algebraic manipulation of this kind calls for high skill and long experience. When a function fails to yield to treatment of this kind it can be often integrated by expanding it into an infinite series, for, as the differential coefficient of a function so expanded is equal to the sum of the differential coefficients of the terms of the series, so the integral of an expanded function can often be obtained by integrating the equivalent series, term by term.

As the differential coefficients of x^n is nx^{n-1} so the integral of nx^{n-1} is x^n . By a little thought we can soon discover what the integral of x^n must be; it will be a power of x , 1 greater than n , for differentiation will decrease this power from $n+1$ to n , and as differentiation will give multiplication by this power, it must be $\frac{x^{n+1}}{n+1}$; for differentiating this we get

$\frac{n+1}{n+1} x^{n+1-1} = x^n$. Thus we easily arrive at a formula for

integrating x^n . Symbolically $\int x^n dx = \frac{x^{n+1}}{n+1} + C$. Suppose

$x=0$ so that $y=1$, the integral is $\frac{x^{0+1}}{0+1} = x$, and the integral of 1 is equal to x . Generally the integral of a constant C is Cx . The reader ought to interpret this geometrically; if $y=C$, a constant, what is the area function of y ?

Now suppose $x=-1$, the integral is $\frac{x^{-1+1}}{-1+1} = \frac{1}{0}$; an

unintelligible result. We might have expected this, for we have already noticed that no power of x has a differential coefficient of x , and we have had to define a function $y = \log x$ to fill this gap. We can now put this definition as $\int \frac{dx}{x} = \log x$, and the reader who remembers our discussion of the graph of the reciprocal function $y = x^{-1}$ on p. 117 will see that what we have called $\log x$ is the area function of a particular rectangular hyperbola.

From our definition of $y = \exp x$ by the relation $\frac{dy}{dx} = y$ we know at once, of course, that the integral of $\exp x$ is simply $\exp x$. The reader ought, using the formula for the integration of x^n , to check that the series purporting to represent this function is reproduced, when it is integrated term by term. If he carries out this calculation he will find that all the series is reproduced excepting the initial term 1. This can be considered to be included in the constant term which, as we have seen, may occur in all indefinite integrals.

The integration of simple harmonic functions gives no difficulties, for as $\frac{dy}{dt}$, for $y = a \sin \omega t$ is $a\omega \cos \omega t$, we see at once that the integral of $a \cos \omega t$ must be $\frac{a}{\omega} \sin \omega t$; similarly the integral of $a \sin \omega t$ must be $-\frac{a}{\omega} \cos \omega t$. Geometrically, the integral of a harmonic function, generated by an angular velocity ω , is obtained by dividing the complex number representing the vector of this function by $j\omega$.

The integration of algebraic functions containing square roots, like $\sqrt{ax^2 + bx + c}$ or fractions like $\frac{ax + b}{ox^2 + dx + e}$, can only be carried out by the kind of artifices of which we made mention. It is outside the scope of this book to deal at length with the technique of integration, and full information on this branch of mathematics can be found in works on the calculus. There are two interesting points that, in this connection, may be briefly explained here. Consider the function $\log f(x)$. The differential coefficient of this is easily

obtained by the rule for a function of a function. We differentiate $\log f(x)$ with respect to $f(x)$ and obtain $\frac{1}{f(x)}$, and multiply this by the differential coefficient of $f(x)$ with respect to x which we may call $f'(x)$. The answer is thus $\frac{f'(x)}{f(x)}$. This answer is very important, as it shows that if an algebraic function, or, for the matter of that, any function is, or can be manipulated, so that it is a fraction with the numerator the differential coefficient of the denominator, then the integral is simply, \log denominator. Thus $\frac{a}{x+b} = a \times \frac{1}{x+b}$ and 1 is the differential coefficient of $x+b$, so that the integral of $\frac{a}{x+b} = a \log (x+b) + C$. Again, $\frac{x}{a^2+x^2} = \frac{1}{2} \frac{2x}{a^2+x^2}$, $2x$ is the differential coefficient of a^2+x^2 , so the required integral is $\frac{1}{2} \log (a^2+x^2) + C$. Algebraic fractions of quite a complicated nature can often, by a special technique of algebra, be decomposed into the sum of what are called partial fractions which satisfy the condition that the differential coefficients of the denominator is equal to the numerator. It is for this reason that, as the reader will find if he pursues his studies beyond the scope of this book, the integration of algebraic functions often leads to logarithmic functions.

Again, we have seen that the differential coefficient of $\arcsin \frac{x}{a}$ and $\arccos \frac{x}{a}$ is $\frac{1}{\sqrt{(a^2-x^2)}}$ and the differential coefficient of $\arctan x$ is $\frac{a}{a^2+x^2}$. These algebraic functions which are differential coefficients therefore give inverse circular functions when they are integrated. This is the basis of a second kind of artifice for the integration of algebraic functions, manipulating or decomposing them so that they correspond to these differential coefficients and yield an inverse circular function as all or part of the integral, and it will account for what the reader will often find in mathematical work, that the integral of an algebraic function containing expressions like $\sqrt{(a^2-x^2)}$ or $\frac{a}{a^2+x^2}$ contains or consists of inverse circular

functions. As a matter of fact it can be seen on reflection that an expression like $\sqrt{(a^2 - x^2)}$ or $(a^2 + x^2)$ is geometrically connected with a variable right-angled triangle containing a constant side a , and if the reader will refer back to p. 124 he will see that we deduced, on geometrical grounds, an area function or integral of an algebraic function that contained a term which was an inverse circular function.

Applications of Integration. The most obvious application of the integral calculus is, of course, the determination of the areas bounded by curves which are the graphical representations of mathematical functions. As an example of this,

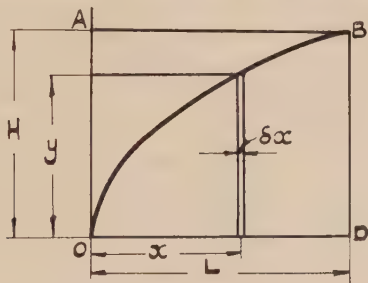


FIG. 68.

consider Fig. 68, which shows the graph of half of the parabolic curve, $x = \frac{y^2}{a}$ or $y = \sqrt{ax}$, which is contained in a rectangle of length L and height H . The angular point B of the rectangle is a point on the parabola and fits its equation, and it has L for its x , and H for its y co-ordinate so that $H = \sqrt{aL}$ and $a = \frac{H^2}{L}$. The area under the curve, being made up of strips like that shown, of height y and width δx , is the definite integral $\int_0^L \sqrt{a \times x} dx$. The indefinite integral is, by the x^n rule, $\frac{2}{3} \times a^{\frac{1}{2}} \times x^{\frac{3}{2}}$; this is 0 for $x = 0$ and $\frac{2}{3} \times a^{\frac{1}{2}} \times L^{\frac{3}{2}}$ for $x = L$. By substituting for a its value in terms of H and L we find that the area OBD contained by the curve, and the straight lines OD and DB is $\frac{2}{3} \times \frac{H}{L^{\frac{1}{2}}} \times L^{\frac{3}{2}} = \frac{2}{3} H \times L$, or $\frac{2}{3}$ of that of the rectangle of length L and height H .

If the figure OBD be conceived to rotate about the line OD it will sweep out a volume known as a paraboloid of revolution, and this volume will be the sum of the volumes swept out by the strips $y \times \delta x$ of which the figure is conceived to be composed. Each of these slices will have a volume of $2\pi y^2 \times \delta x$, for they are circular, of radius y , and thickness δx . The total volume will be given by the integral $\int_0^L 2\pi y^2 dx$ or $\int_0^L 2\pi \times \frac{H^2}{L} \times x dx$, or $2\pi \times \frac{H^2}{L} \times \frac{1}{2}L^2 = \frac{1}{2} \times 2\pi H^2 \times L$ which is easily seen to be one half of the volume of the cylinder in which the paraboloid is contained.

The notion of the area bounded by a curve leads to the idea of an average height or ordinate of this curve. Thus, in Fig. 68, the area under the curve is equal to a rectangle of length L and height $\frac{2}{3}H$, so the average height of the curve OB is equal to $\frac{2}{3}H$. We have already used this idea on p. 177, in estimating the average velocity of a point executing simple harmonic motion, for, as the reader will see on reflection, the double amplitude or total travel of the point is given by the definite integral of the velocity function.

Another important application of the integral calculus is the determination of what are called first and second moments. Consider Fig. 68 again, the quantity $y\delta x \times x$ is called the first moment of the area of the strip $y\delta x$ about the axis OA, and the integral $\int_0^L yx dx$ is called the first moment of the whole area OBD. We can easily evaluate this integral; it is $\int_0^L \frac{H}{L} \times x^2 dx$, and its value is $\frac{2}{5} \frac{H}{L} \times L^3 = \frac{2}{5}HL^2$.

The first moment of an area, divided by the measure of its area, gives the distance of what is called the centre of area or centroid from the axis about which the first moment was calculated. The distance of the centroid of OBD of Fig. 68 from the vertical OA is therefore $\frac{2}{5}HL^2 \div \frac{2}{3}HL = \frac{3}{5}L$.

The quantity $y\delta x \times x^2$ is called the second moment of the strip of Fig. 68 defined by y and δx , and the integral of this quantity between the limits O and L is called the second moment of the area OBD. This second moment integral is

easily seen to be $\int_0^L \frac{H}{L^3} \times x^3 dx$ and its value is $\frac{2}{7}HL^3$ or $\frac{3}{7}L^2 \times$ area of OBD. If the area OAB is a material and homogenous lamina having a mass m , the quantity m is proportional to the area, $\frac{3}{7}L^2 \times m$ is called the moment of inertia of the lamina, and $\frac{3}{7}L^2$ is the square of the radius of gyration about the axis OA. This means that if the lamina is rotated about the axis OA with uniform angular velocity its kinetic energy will be the same as that of a particle of mass m , distant from OA by $\sqrt{\frac{3}{7}} \times L$, and rotated about OA with the same angular velocity.

A large number of examples of the integral calculus to the establishment of mensurational formulae for areas and volumes and for the determination of first and second moments will be found in any textbook on the calculus. As the discovery of the indefinite integral corresponding to a given function calls for mathematical acumen of a high order, so the application of the formulae for known integrals to specific problems requires considerable care and practice in the algebraic working. As this book is concerned with general principles rather than with the details of mathematical technique, we shall say no more on this matter.

Differential Equations. An implicit function of x , written symbolically $f(x,y)=0$, is also an equation, the solution of which gives y as an explicit function of x . If such an equation contains differential coefficients of y with respect to x it is then called a differential equation, and the function $y=\psi(x)$, that is a solution of or that satisfies the equation is called its integral.

A simple example of a differential equation is $\frac{dy}{dx} - x = 0$.

This is equivalent to $\frac{dy}{dx} = x$ and is simply a statement of the problem; find the integral of x with respect to x . The answer, as we know, is $y = \frac{1}{2}x^2 + C$, because the differential coefficient of this function is x . Again, consider the differential equation $\frac{dy}{dx} - y = 0$. This is equivalent to $\frac{dy}{dx} = y$, and we know, by the definition of $\exp x$, that the equation is satisfied by

$y = \exp x$. But it is also satisfied by $y = A \times \exp x$ where A is a constant, so this more general function is the most complete form of the solution or integral.

Let us take a more difficult example, $\frac{d^2x}{dt^2} + p^2x = 0$. This is $\frac{d^2x}{dt^2} = -p^2x$, and x must be a function of t such that when it is differentiated twice, x is reproduced with $-p^2$ as a multiplier. We know that the simple harmonic function $x = a \sin pt$ satisfies this condition, and we might think that this can be called the solution or integral. But $x = b \cos pt$ satisfies the equation also, and if the mathematical operation denoted by $\frac{d^2x}{dt^2} + p^2x$ when carried out on each of these functions gives 0 as an answer, the answer will still be 0 if it is carried out on their sum. The proper statement of the integral is therefore $x = a \sin pt + b \cos pt$, which, as we saw on p.197, is equivalent to $x = A \sin (pt + \theta)$ where A and θ are defined by a and b in the manner there shown. This is worth a little thought; it means that the differential equation is satisfied by a function containing constants that are quite arbitrary; they can be given any value we choose. An equation, like that which we have just considered, contains a second differential coefficient, and it is said to be of the second order. Its integral contains two arbitrary constants, these are not implicitly contained in the equation because in its formation they disappear, or are eliminated by differentiation. A differential equation is therefore more general in its implication than an ordinary equation. The solution of an ordinary equation gives the unknown as a function of quantities presumed to be known and specified; a differential equation of the n th order gives a solution or integral which contains n new constants, the values of which are quite arbitrary and do not affect the correctness of the solution. If the reader pursues his study of mathematics beyond the scope of this book he will find that the integral of a higher kind of differential equation, with two independent variables, may be still more general; the kind of functions in the integral may be undefined. The differential equation to a wave motion of one dimension is written

$\frac{d^2y}{dx^2} = a^2 \frac{d^2y}{dt^2}$ and the solution or integral of this equation is $y = f(x + at) + \psi(x - at)$, in which the very forms of the functions of $(x + at)$ and $(x - at)$ are arbitrary. We mention this merely as an example of the fact that the higher the class to which a differential equation belongs the more general is the character of its solution or integral; comment or explanation of the equation is outside the scope of this book.

Let us consider, as a final example, the differential equation $\frac{d^2x}{dt^2} + p^2x = B \sin \omega t$. Suppose there is some function of t , $x = f(t)$ such that its second differential coefficient plus p^2 times itself equals $B \sin \omega t$. This satisfies the equation. But we know that if the right-hand side of the equation is 0 it will be satisfied by $x = A \sin(\omega t + \theta)$. It follows therefore that if the calculation indicated by the differential equation be carried out on $A \sin(\omega t + \theta) + f(t)$ the result will be $0 + B \sin \omega t$. The complete solution therefore consists of two parts, the first, called the complementary function, is the integral for $B \sin \omega t = 0$, the second, called the particular integral, is the function which, when subjected, as it were, to the operation indicated by $\frac{d^2x}{dt^2} + p^2x$ gives $B \sin \omega t$ for an answer. We can easily discover this particular integral; it must be a circular function, corresponding to an angular velocity p , to give such a function as the result of the $\frac{d^2x}{dt^2} + p^2x$ operation, and we know that the effect of the operation $\frac{d^2x}{dt^2}$ on a circular function is to multiply it by $-\omega^2$. Thus the operation $\frac{d^2x}{dt^2} + p^2x$ is the same as $(p^2 - \omega^2)x$, so the differential equation, as far as the particular integral is concerned, is simply $(p^2 - \omega^2)x = B \sin \omega t$, so that $x = \frac{B}{p^2 - \omega^2} \sin \omega t$. The complete solution is therefore $x = A \sin(pt + \theta) + \frac{B}{p^2 - \omega^2} \sin \omega t$.

The differential equations encountered in applied mathematics are symbolic methods of writing quantitative state-

ments about natural phenomena controlled by physical laws.

Thus, as we saw some time back, an equation like $m \frac{d^2x}{dt^2} + Fx = 0$,

where m stands for a mass, and F for a control force per unit displacement, is one all terms of which have the physical dimensions of a force, and it states that acceleration forces are at any instant exactly balanced by control forces. The integral of this equation shows that the motion represented by x is simple harmonic, and that its period depends upon F and

m . If this equation is altered to $m \frac{d^2x}{dt^2} + Fx = B \sin \omega t$ this

shows that in addition to acceleration and control forces we have a third variable force represented by $B \sin \omega t$, the period of which is different from that depending upon m and F . m and F define a free period corresponding to the absence of this third force; the forced period corresponding to ω , tends to change the simple harmonic motion defined by m and F . We can see easily that the free period is connected with the angular velocity p , in the solution of the equation given in the preceding paragraph, and this solution shows that when ω becomes nearly equal to p , or, when the forced and free periods are nearly

alike, the amplitude $\frac{B}{p^2 - \omega^2}$ of the particular integral term becomes very large and we have the mechanical phenomenon of resonance.

The foregoing is but a slight sketch written with a view to giving the reader an idea of what a differential equation is. The solution of differential equations calls for mathematical skill of an order even higher than that required for integration, and it forms the subject of a well-defined branch of advanced mathematics.

CHAPTER X

EXPONENTIAL AND LOGARITHMIC FUNCTIONS

Area Function of a Rectangular Hyperbola. We have defined two functions, $\log x$ and $\exp x$ by their differential properties; the differential coefficient of $\log x$ is $\frac{1}{x}$, that of $\exp x$ is $\exp x$. The geometrical meaning of this definition of $\log x$ is that it is the area function of the graph of the function $y = \frac{1}{x}$, which we met with in Chapter VI and which is called a rectangular hyperbola. We have now to make a somewhat detailed study of this area function which is the geometrical representation of what we have called $\log x$.

In Fig. 69 we have the graph of the rectangular hyperbola, the equation to which is $y = \frac{1}{x}$. This is symmetrical with respect to the 45 degree line bisecting the angle between the axes of co-ordinates, and it has the geometrical property that the area of any rectangle defined by the co-ordinates of any point is equal to the area of the square OCDB, which is unity. Thus the rectangles OP_2 , and OP_1 are each of unit area. Now we know that the integral of a function like $\frac{1}{x}$ is indefinite in its value, since it includes an arbitrary constant, therefore the measure of the area function of the corresponding graph must be defined in reference to some fixed vertical ordinate. We shall define the area function corresponding to a point such as P_1 as the area contained between the ordinate LP_1 , the horizontal axis, the ordinate at B where $x=1$, and the graph. This area function is numerically equal to the measure of the full-line shaded area in the figure. Thus if x or OL is measured in inches, the number of square inches contained in this shaded area is the value of the area function of x . We can denote this by $\log x$, but we shall use the notation $A(x)$ for this function to denote that we are con-

sidering it from a geometrical standpoint. It follows from this definition that when $x=1$, $A(x)=0$, and $\log 1=0$.

If $x=2$, the right-hand extremity of the associated area will be $\frac{1}{2}$ and the area function will be greater than the area of a rectangle $1 \times \frac{1}{2}$ or $\frac{1}{2}$ units of area. If x is increased to 3, the right-hand extremity of the area function will be $\frac{1}{3}$ and the area function will be increased by an amount greater

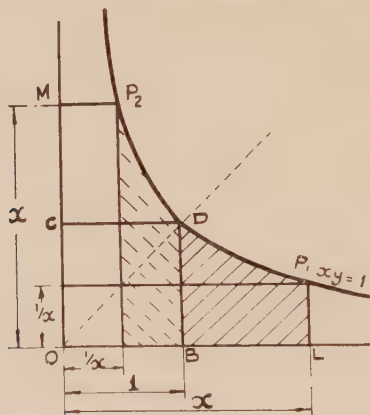


FIG. 69.

than $\frac{1}{3}$. Thus if $x=n$ the area function will be greater than the sum of the series $\frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} \dots \frac{1}{n}$. We have seen on p. 103 that this sum increases without limit as n is increased indefinitely, and we conclude that *a fortiori*, $A(x)$ or $\log x$ increases without limit as x so increases. This should not be taken as obvious. The graph of $\frac{1}{x^2}$ looks much like that of $\frac{1}{x}$, but the area function of $\frac{1}{x^2}$ can be shown to approach the limit 1, as x increases indefinitely. The increase of $A(x)$ with x is geometrically continuous, hence $A(x)$ or $\log x$ is a continuous function of x .

Let OM in Fig. 69 be equal to OL; M defines a point P_2 on the graph, the horizontal ordinate of which is $\frac{1}{x}$. The

vertical ordinate through P_2 defines an area function, shown dotted-line shaded, which corresponds to $A\left(\frac{1}{x}\right)$. This area lies to the left of BD , and it is reckoned negative. If $y = A\left(\frac{1}{x}\right) = \log\left(\frac{1}{x}\right)$ then, by the rules for the differentiation of a function of a function, $\frac{dy}{dx} = \left(1 \div \frac{1}{x}\right) \times -\frac{1}{x^2} = -\frac{x}{x^2} = -\frac{1}{x}$ = minus the differential coefficient of $-\log x$. We conclude, therefore, that $\log \frac{1}{x} = -\log x$.

The reader will readily see, by a little careful study of Fig. 69, that $A\left(\frac{1}{x}\right)$ is equal numerically to $A(x)$, for the figure CMP_2D is manifestly congruent to $A(x)$, full line dotted, and the rectangle OP_2 is equivalent to the square $OCDB$. It follows, therefore, that as $\frac{1}{x}$ gets smaller and smaller by increase of x , $A\left(\frac{1}{x}\right)$ or $\log \frac{1}{x}$ increases numerically without limit in a negative direction. The area function of zero, or $\log 0$ is an indefinitely large negative number. We can, at present, give no meaning for the area function of a negative quantity.

Now, consider Fig. 70. Here is the graph of $y = \frac{1}{x}$ plotted to a smaller scale than it is in Fig. 69. The area function corresponding to x_1 is shown divided into a number of strips all of equal width, say a . A point on the horizontal axis defined by the length x_2 , will have a corresponding vertical ordinate equal to $\frac{1}{x_2}$. Suppose a strip be marked off to the right of this ordinate of width $a \times x_2$. If a is very small the area of the first strip of the area function of x_1 will be a , for its height is 1. The area of the new strip at x_2 will also be a . The second strip of the area function of x_1 will have a starting height of $\frac{1}{1+a}$ and its area will be $\frac{a}{1+a}$. If a second strip of width $a \times x_2$ be

Raising this number to the n th power gives $(x^{\frac{1}{n}})^n = x$. Thus $A\left[\left(x^{\frac{1}{n}}\right)^n\right] = Ax = nAx^{\frac{1}{n}}$, so that $A(x^{\frac{1}{n}}) = \frac{1}{n}A(x)$. We can therefore conclude that $A(x^{\frac{m}{n}}) = \frac{m}{n}A(x)$, so that when n is any positive rational number $A(x^n) = nA(x)$. Further, as $A\left(\frac{1}{x^n}\right) = -A(x^n)$, this rule holds, whether n is positive or negative. Finally, as $A(x)$ is a continuous function of x we may conclude then $nA(x) = A(x^n)$ when n is any positive real number, or rather,

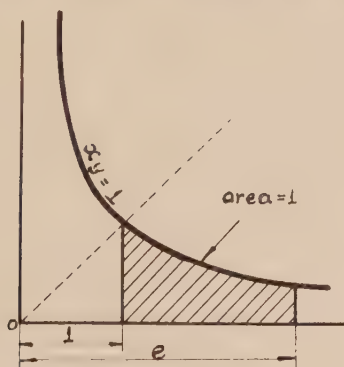


FIG. 71.

that the meaning of x^n is that its area function is n times that of x .

Suppose that $A(x) = 1$. This will define a numerical value of x , which is denoted in mathematics by the symbol e . The geometrical significance of the number is shown in Fig. 71. It is the value of the x ordinate, measured in inches, that gives an area function of exactly 1 square inch.

We have used the symbol $A(x)$ as a mathematical synonym for what was defined originally as $\log x$, and all we have proved for $A(x)$ is, of course, true for $\log x$. e is therefore a number such that $\log e = 1$. Thus $\log x = \log x \times \log e$ or $\log e^{\log x}$, so that $x = e^{\log x}$.

The function $y = \exp x$ has been defined by the equation

$\frac{dy}{dx} = y$, or $\frac{dx}{dy} = \frac{1}{y}$, so that x is the logarithm of y . It follows therefore from the last paragraph that $y = e^{\log y} = e^x$, so that the $\exp x$ function is equivalent to e^x , where x has all the attributes of an index or power.

Exponential Theorem. We have now elucidated fairly completely the nature of the functions we have defined as $\log x$ and $\exp x$, and we had better pause and consider carefully what this all means. We started our consideration of these functions in the last chapter, when we assumed that these functions exist and stipulated their differential properties. The geometrical interpretation of these properties has led us, by a

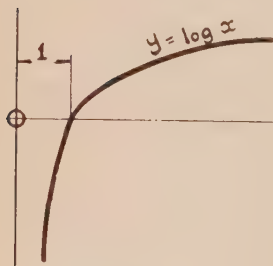


FIG. 72.

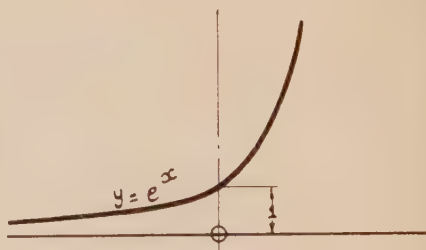


FIG. 73.

study of the rectangular hyperbola to which they are related, to conclude that $\log x$ is a function having all the algebraic properties of the index of a power, while $\exp x$ has all the properties of an exponential function, that is, a constant number e raised to any power x . $\log x$ is the index of the power e that is equal to x , $\log x_1 x_2$ is equal to $\log x_1 + \log x_2$ and $\log x^n$ is equal to $n \log x$, in the same way that $e^a \times e^b = e^{a+b}$ and $(e^a)^n = e^{an}$ by the elementary index rule of algebra. Objectively $\log x$ is the measure of the area function corresponding to a horizontal ordinate of the curve represented by $xy=1$, while e is the value of this ordinate that has an area function equal to a . The graph of $y = \log x$ is as shown in Fig. 72; as x approaches zero, y becomes indefinitely large and negative. There is no real value of y when x is negative. The graph of $y = e^x$ is as shown in Fig. 73; as x increases without

limit in a negative direction y tends to zero. The horizontal axis is an asymptote to the graph.

In the preceding chapter, on p. 201, we showed that the differential property of the $\exp x$ function enabled us to expand the function into an infinite convergent series. Since $\exp x$ is equivalent to e^x we can now write :

$$\exp x = e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!}.$$

This is called the exponential theorem. It implies that if we evaluate two numerical values of the series for, say x_1 and x_2 , the product of these two calculated numbers will be equal to that third number we should obtain by calculating the value of the series for $(x_1 + x_2)$. Similarly, the n th power of the calculated value of the series for any number x will be equal to the value calculated for $(n \times x)$.

The exponential series enables us rapidly to obtain an approximation to the numerical value of e . For

$$e = e^1 = 1 + 1 + \frac{1}{2!} + \frac{1}{3!} + \dots$$

We can set out this calculation as follows :

$$\begin{array}{r} 1 \cdot 000\ 000 \\ 2) 1 \cdot 000\ 000 \\ 3) \cdot 500\ 000 \\ 4) \cdot 166\ 666 \\ 5) \cdot 041\ 666 \\ 6) \cdot 008\ 333 \\ 7) \cdot 001\ 389 \\ 8) \cdot 000\ 198 \\ 9) \cdot 000\ 025 \\ \quad \cdot 000\ 003 \\ \hline e = \quad 2 \cdot 718\ 280 \end{array}$$

The square of this number, or e^2 , ought to be obtained by using 2 for x in the e^x series. We made this calculation on p. 104, and the reader will not be wasting his time if he tests by actual arithmetic the correspondence between the square of 2.718 280 and the series value of e^2 .

Logarithms. The importance of the logarithms' function, as is well known, is that it is the basis of a method of arithmetical calculation in which addition, subtraction, multiplication, and division are respectively substituted for multiplication, division, evolution, and involution. Thus if the values of $\log x_1$ and x_2 are known $x_1 + x_2$ is the logarithm of $x_1 \times x_2$. $x_1 \times x_2$ is called the antilog function of $x_1 + x_2$, and if complete tables of the log and antilog functions are available, the product $x_1 \times x_2$ can be evaluated by the mere addition of logarithms. The technique of this calculation is explained in all books on practical mathematics.

The value of the function $y = \log x$ is known as the logarithm of x as defined above, that is, as the index of a power of e . This kind of logarithm is sometimes called the natural, the hyperbolic, or the Napierian logarithm. In mathematics, however, $\log x$ always implies a power of e , without any qualifying description. Practical calculations are, for reasons explained on p. 68, made with what are called common logarithms, which are equivalent to the indices of powers of 10. The common logarithm of x is denoted by the symbol $\log_{10} x$. There is a very simple connection of the two classes of logarithms. If $a = \log 10$ $10 = e^a$. The common logarithm b of x is defined by $x = 10^b = (e^a)^b = e^{ba}$, so that $\log x = a \times b$, or b the common logarithm $= \log x$ divided by $\log 10$. We may mention here, without proof and merely as a matter of interest, that logarithms to the base 10 are area functions of a hyperbola of which the asymptotes are inclined at an angle less than 90 degrees.

It is evidently impossible to express the function $\log x$ as an infinite series of the form, $a + bx + cx^2 \dots$, for no series of this form can give an infinite negative value of the function when $x = 0$. We can, however, expand the function $\log (1 + x)$ which becomes $\log 1 = 0$ when $x = 0$, subject to certain restrictions. The series $1 - x + x^2 - x^3 \dots$ is a geometrical progress, the common ratio of which is $-x$, and, according to the rule given on p. 102 the sum to infinity of this series, provided x is numerically less than 1, is $\frac{1}{1+x}$. Thus $\frac{1}{1+x} = 1 - x + x^2 - x^3 \dots$ subject to this condition. Now, as the

differential coefficient of the fraction $\frac{1}{1+x}$ is equal to the denominator, the integral of the fraction is $\log(1+x)$, and if we can assume that this integral is equal to the integration of the series, term by term, we have, as the integral of 1 is x , that of x is $\frac{x^2}{2}$, and so on :

$$\log(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4}.$$

But, we ought to ask, what about the unknown constant of integration? This must be such that the series satisfies the fundamental condition $\log 1 = 0$, and we see at once that no constant is required. The series is true for values of x less than 1, and it can be proved by higher algebra to be true for $x = 1$ provided the order of the terms of the series is unchanged. If $x = 1$ the series gives $\log 2$, and, if the reader will refer back to p. 103, we will see that this $\log 2$ series is identical with the one for which we drew attention to some very curious properties.

The series for $\log(1+x)$ converges very slowly and is unsuited for the numerical calculation of $\log x$. It is shown in books on algebra how this series can be used to form others which are of greater practical utility, and the reader will also find there, if he is interested, that there are other methods of calculating logarithms.

Hyperbolic Functions. The circular functions, $\sin \theta$, $\cos \theta$ and $\tan \theta$, are defined geometrically in reference to an independent variable θ which is usually stated as the ratio of the arc of a circle to its radius a , but which, as we pointed out on p. 48, can also be considered to be connected with the sectorial area A , contained by it according to the relation $A = \frac{1}{2}a^2\theta$. We have seen that the curve whose equation is $y = \sqrt{(x^2 - a^2)}$ can, in a way, be considered to be a sort of extension of the circle corresponding to $y = \sqrt{(a^2 - x^2)}$ when x takes values greater than the radius a . The $y = \sqrt{(x^2 - a^2)}$ curve was stated on p. 117 to be a rectangular hypobola, and shown to be identical with the curve represented by $xy = \frac{a^2}{2}$,

twisted through 45 degrees. We have now to consider an important class of functions, which are related to the hyperbola of $y = \sqrt{(x^2 - a^2)}$ in a manner analogous to the relation of the sin and cos functions to the circle of $y = \sqrt{(a^2 - x^2)}$.

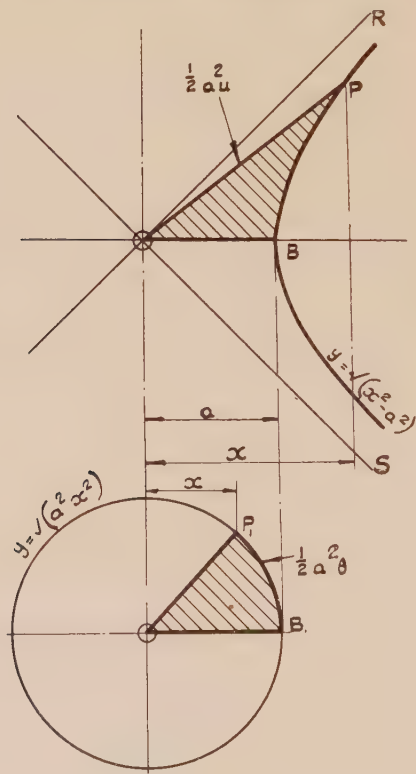


FIG. 74.

The upper graph of Fig. 74 is one branch of the rectangular hyperbola whose equation is $y = \sqrt{(x^2 - a^2)}$. The point P on the curve defines x and y co-ordinates, it also defines a sectorial area, shown shaded, bounded by the curve, the line OB , and the radial line OP . This area is analogous to the sectorial area associated with a geometrical angle in a circle.

If the sectorial area A has a measure of $\frac{1}{2}a^2u$, so that $u = \frac{2A}{a^2}$,

then the ratio $\frac{x}{a}$ corresponding to the point P is called the hyperbola cosine of the quantity u defined by the sectorial area associated with P . This is written $\cosh u = \frac{x}{a}$. $\frac{y}{a}$ is called the hyperbolic sine of u and is written $\sinh u = \frac{y}{a}$. The ratio

$\frac{y}{x} = \frac{\sinh u}{\cosh u} = \tanh u$. These \cosh , \sinh , and \tanh ratios are functions of an independent variable u , which evidently stands in the same relation to the hyperbola as the θ in $\cos \theta$ stands to a circle. We have to elucidate the nature of these so-called hyperbolic functions.

Fig. 75 is the graph of the function $y = \sqrt{x^2 - 2}$, for which the a length, OB , is $\sqrt{2}$, and we have seen that the graph also corresponds to $xy = 1$ if the asymptotes OR and OS are the axes of co-ordinates. The full-line shaded area is, according to our definition of u , $\frac{(\sqrt{2})^2 u}{2} = u$. As $\cosh u$ is defined as

$\frac{x}{a} = \frac{x}{\sqrt{2}}$, the x co-ordinate corresponding to the area u is $\sqrt{2} \cosh u$. Similarly the y ordinate is $\sqrt{2} \sinh u$. Now the point P defines x and y co-ordinates of the graph relative to the axes OR and OS . This x co-ordinate is OL and a little study of the diagram will show that as OL is inclined 45 degrees to OX , the distance OL is equal to $\frac{1}{\sqrt{2}}$ the distance

OX or $\frac{x}{\sqrt{2}}$ less $\frac{1}{\sqrt{2}}$ the distance XP or $\frac{y}{\sqrt{2}}$. Thus $OL = \frac{x}{\sqrt{2}} - \frac{y}{\sqrt{2}}$ or $\cosh u - \sinh u$. Now, considering the curve in relation to the asymptote axes OR and OS , the area function corresponding to OL , or $\log(OL)$ is the area shown dotted-line shaded and a little further study of the diagram will show that this area, $\log(OL)$, or $\log(\cosh u - \sinh u)$ is equal to the full-line shaded area u , for, of the parts of the areas that do not overlap, the triangle OPL is equivalent to the half-square OBM , by the $xy = 1$ property. But the dotted-line shaded

area is reckoned negative, so that $\log (\cosh u - \sinh u) = -u$, and $\cosh u - \sinh u = e^{-u}$. Further, as $y^2 = x^2 - 2$, $x^2 - y^2 = 2$, $2 \cosh^2 u - 2 \sinh^2 u = 2$, and $\cosh^2 u - \sinh^2 u = 1$, which can be converted to $(\cosh u - \sinh u)(\cosh u + \sinh u) = 1$. Thus

$$\cosh u + \sinh u = \frac{1}{\cosh u - \sinh u} = \frac{1}{e^{-u}} = e^u. \text{ We thus have,}$$

$$\cosh u - \sinh u = e^{-u} \text{ and } \cosh u + \sinh u = e^u.$$

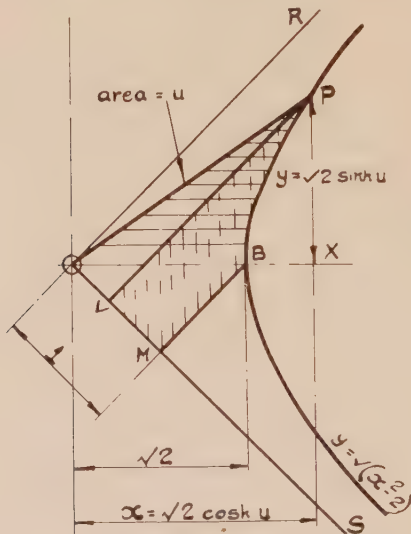


FIG. 75.

Whence, adding and subtracting the equations :

$$\cosh u = \frac{1}{2}(e^u + e^{-u}), \text{ and } \sinh u = \frac{1}{2}(e^u - e^{-u}).$$

It is easy to deduce some properties of the hyperbola functions from their fundamental definitions given above; thus when $u=0$ $\cosh u=1$ and $\sinh u=0$. As u increases in Fig. 75, it is evident that $\cosh u$ and $\sinh u$ increase without limit. As P approaches the asymptote for very large distances from O, the ratio $\frac{\sinh u}{\cosh u} = \tanh u$, tends to the limiting value 1, for the x and y co-ordinates of the graph tend to equality.

As we have found the series equivalent to e^u , and can write the series for e^{-u} by substitution of $(-u)$ for u , we have

$$e^u = 1 + u + \frac{u^2}{2!} + \frac{u^3}{3!} + \frac{u^4}{4!} \dots$$

and

$$e^{-u} = 1 - u + \frac{u^2}{2!} - \frac{u^3}{3!} + \frac{u^4}{4!} \dots$$

Adding the two series, odd powers disappear, and even powers are doubled so that

$$\cosh u = \frac{1}{2}(e^u + e^{-u}) = 1 + \frac{u^2}{2!} + \frac{u^4}{4!} \dots$$

Subtracting the second from the first series, even powers disappear, and odd powers are doubled and

$$\sinh u = \frac{1}{2}(e^u - e^{-u}) = u + \frac{u^3}{3!} + \frac{u^5}{5!} \dots$$

These series show that $\cosh u$ is an even, and $\sinh u$ an odd

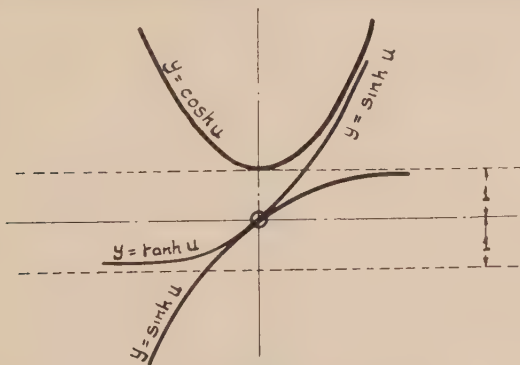


FIG. 76.

function of u . The graphs of the three hyperbolic functions are shown in Fig. 77.

The differential properties of these functions are easily obtained. Thus if $y = \cosh u = \frac{1}{2}e^u + \frac{1}{2}e^{-u}$, $\frac{dy}{du} = \frac{1}{2}e^u - \frac{1}{2}e^{-u} = \sinh u$, and $\frac{d^2y}{du^2} = \frac{1}{2}e^u + \frac{1}{2}e^{-u} = \cosh u = y$. Similarly we find that if $y = \sinh u$, $\frac{dy}{du} = \cosh u$, and $\frac{d^2y}{du^2} = \sinh u = y$. Thus $\cosh u$

and $\sinh u$ are functions which are reproduced exactly by two differentiations.

Inverse Hyperbolic Functions. If $x = \cosh u$, u is an inverse hyperbolic functions of x which is written $u = \cosh^{-1} x$. A similar meaning applies to $u = \sinh^{-1} x$, and $\tanh^{-1} x$. The first two functions have a real value for all positive values of x , but $\tanh^{-1} u = x$ has a meaning only when x is not greater than 1.

If $x = \cosh u$, then $\frac{dx}{du} = \sinh u$, and, as $\cosh^2 u - \sinh^2 u = 1$
 $\sinh u = \sqrt{(\cosh^2 u - 1)} = \sqrt{(x^2 - 1)}$. Therefore $\frac{du}{dx}$, or the differential coefficient of $u = \cosh^{-1} x$ is $\frac{1}{\sqrt{(x^2 - 1)}}$. We find similarly that the differential coefficients of $\sinh^{-1} x$ and $\tanh^{-1} x$ are respectively $\frac{1}{\sqrt{(x^2 + 1)}}$ and $\frac{1}{1 - x^2}$. These derived functions of the inverse hyperbolic functions are algebraic, and because of this the integrals of many algebraic functions can be expressed all or in part as inverse hyperbolic functions.

These functions can themselves be converted to the logarithmic form. Since $u = \log (\cosh u + \sinh u)$ and $\sinh u = \sqrt{(\cosh^2 u - 1)}$ it follows that if $\cosh u = x$, then $u = \cosh^{-1} x = \log [x + \sqrt{(x^2 - 1)}]$. Similarly if $u = \sinh^{-1} x$, then u also equals $\log [x + \sqrt{(x^2 + 1)}]$. Again, $u = \log (\cosh u + \sinh u)$ and it also equals $-\log (\cosh u - \sinh u)$, and adding these two equations we find that $2u = \log (\cosh u + \sinh u) - \log (\cosh u - \sinh u)$ or $\log \frac{\cosh u + \sinh u}{\cosh u - \sinh u}$. If we divide

every term in the fraction by $\cosh u$ and put $\tanh u = \frac{\sinh u}{\cosh u} = x$,

then if $u = \tanh^{-1} x$ we find that $u = \frac{1}{2} \log \frac{1+x}{1-x}$. Because of

these alternative forms of the inverse hyperbolic functions the reader will find that, in the textbooks, the integrals of the same algebraic functions are given sometimes as logarithmic and sometimes as $\cosh^{-1} x$, $\sinh^{-1} x$, or $\tanh^{-1} x$ functions.

Exponential Expressions for the Circular Functions. Let us now turn back to Fig. 74, and consider the changes of the area function of $y = \sqrt{(x^2 - a^2)}$. When $x = a$, the radial line OP coincides with OB and the area function is zero. If x is less than $a = OB$, we have seen that the values of y corresponding become imaginary, and, in the range of x values from $+a$ to $-a$, the corresponding y values correspond to those defined by a circle of radius a lying in a perpendicular plane. This circle is shown in the lower part of the figure which may be considered to be a plan of the upper part. When x is less than a it will define an area function of the circle which is shown shaded, and if $a = \sqrt{2}$, as it is in Fig. 76, the numerical value of the circle area function will be $\frac{1}{2} \times (\sqrt{2})^2 \times \theta$, or simply θ , where θ is the radian measure of the angle defined by the x value. Now when x is greater than $\sqrt{2}$ the area function of the hyperbolic curve has been shown to be the logarithm of $\cosh u + \sinh u$, or of $\frac{x}{a} + \frac{y}{a}$.

Can we make any similar statement when x is less than $a = \sqrt{2}$, and it defines a point on the imaginary circle? In this condition we see at once that the ratios $\frac{x}{a}$ and $\frac{y}{a}$, called $\cosh u$ and $\sinh u$ for the hyperbola are $\cos \theta$ and $\sin \theta$ for the circle, for which θ corresponds to an imaginary area function of the hyperbola. Let us assume that the relation we found for the hyperbola $y = \sqrt{(x^2 - 2)}$, that $u = \log (\cosh u + \sinh u)$ is, by a kind of analogy, true for the imaginary form of the hyperbola, the circle $y = \sqrt{(2 - x^2)}$, and that we indicate that the circle is imaginary by calling its area function $j\theta$, and the ratio $\frac{y}{a}$ which contains the distance y in the imaginary plane, $j \sin \theta$. If this assumption is justified, then we ought to have $j\theta = \log (\cos \theta + j \sin \theta)$. We are here getting into rather deep waters by giving $\cos \theta + j \sin \theta$, an imaginary logarithm, because this at present has no intelligible meaning. If, however, we ignore this and go ahead, by assuming that an imaginary logarithm has the properties of a real one we can say $\cos \theta + j \sin \theta = e^{j\theta}$, and just as we found exponential expressions for $\cosh u$ and

$\sinh u$, so we can show that, according to the rules of algebra,

$$\cos \theta = \frac{1}{2}(e^{j\theta} + e^{-j\theta}), \text{ and } \sin \theta = \frac{1}{2j}(e^{j\theta} - e^{-j\theta}).$$

These statements are still meaningless algebraically, because raising e to an imaginary power is an unintelligible operation. Let us, however, still go ahead and assume that, as e can be expanded into a series of power of x , so $e^{j\theta}$ can be expanded into a series of powers of $(j\theta)$. We see at once that $(j\theta)^2 = -\theta^2$, $(j\theta)^3 = -j\theta^3$, $(j\theta)^4 = \theta^4$, $(j\theta)^5 = j(j\theta)^4 = j\theta^4$, and so on. Thus, as $-j\theta = j(-\theta)$

$$e^{j\theta} = 1 + j\theta - \frac{\theta^2}{2!} - \frac{j\theta^3}{3!} + \frac{\theta^4}{4!} + \frac{j\theta^5}{5!} \dots$$

$$e^{-j\theta} = 1 - j\theta - \frac{\theta^2}{2!} + \frac{j\theta^3}{3!} + \frac{\theta^4}{4!} - \frac{j\theta^5}{5!} \dots$$

as the signs of even powers of $j\theta$ and $j(-\theta)$ will be like, and the odd powers unlike. Adding the two series, odd powers of θ disappear and even powers are doubled so that

$$\cos \theta = \frac{1}{2}(e^{j\theta} + e^{-j\theta}) = 1 - \frac{\theta^2}{2!} + \frac{\theta^4}{4!} \dots$$

Subtracting the second from the first series even powers disappear, all odd powers are doubled, and these all contain j as a factor. Thus,

$$\sin \theta = \frac{1}{2j}(e^{j\theta} - e^{-j\theta}) = \theta - \frac{\theta^3}{3!} + \frac{\theta^5}{5!} \dots$$

The reader will see that these series, obtained by assuming that $e^{j\theta}$ can be expanded into a power series are identical with those deduced on p. 199 from the differential properties of the circular functions, which properties were themselves deduced earlier on p. 176 from elementary geometrical considerations.

What does this all mean? To answer this question we must retrace our steps and summarise our chain of argument. We have, from the geometry of the hyperbola deduced a logarithmic relation between its area function and its x and y co-ordinates which can be expressed in our symbolism as $u = \log (\cosh u + \sinh u)$. We have considered the circle to

be an imaginary form of the hyperbola, and we have assumed that a similar relation exists between the imaginary area function represented by the θ of the circle, its real x , and its imaginary y co-ordinate. This assumption leads to the result that $j\theta = \log(\cos \theta + j \sin \theta)$ or $\cos \theta + j \sin \theta = e^{-j\theta}$. We have then assumed that the exponential function $e^{j\theta}$, with its imaginary exponent, can be expanded into a power series by the same rule as has been shown to apply when the exponent of e is real. The result of this chain of assumptions is that we have obtained series values for $\cos \theta$ and $\sin \theta$ that are identical with those obtained in quite another way, and without the use of any exponential or logarithmic functions. The answer to our question is, therefore, that as we found that imaginary and complex numbers can be treated algebraically as can ordinary or real numbers, so imaginary exponents and logarithms, if assumed to have all the properties of those that are real, lead to consistent and intelligible results.

The equation $\cos \theta + j \sin \theta = e^{j\theta}$ is a very remarkable one of far-reaching implication. Suppose $\theta = \pi$ (180 degrees) $\sin \theta = 0$ and $\cos \theta = -1$, so that $-1 = e^{j\pi}$. If $\theta = 2\pi$, $\cos \theta = +1$ and $\sin \theta = 0$ so that $e^{2j\pi} = 1$. As we have found a value for $\log -1$, $j\pi$, we can assign a value to the logarithm of any negative number $-N$, for $-N = -1 \times (+N)$ and $\log -N = \log -1 + \log (+N)$ so $\log -N = j\theta + \log N$. We can go further than this, for as $+N = +1 \times N$ we can say that $\log N = \log (+1) + \log N = 2\pi j + \log N$. Indeed, every positive number has an infinite number of logarithms, one, $\log N$, is real, and the others obtained by successive additions of $2j\pi$. The logarithms of a negative number $-N$ can be obtained by successive additions of $2\pi j$ to $j\pi + \log N$.

The equation $\cos \theta + j \sin \theta = e^{j\theta}$ also shows that our tentative ideas of p. 166 are justified. For, if the θ of $\cos \theta + j \sin \theta = e^{j\theta}$ is increased from 0, when $e^{j\theta} = e^0 = 1$, a unit number in the real line, the effect of this increase is to alter the inclination of the line represented by the complex number $\cos \theta + j \sin \theta$, without altering its length. Raising e to an imaginary power, therefore, is simply equivalent, geometrically, to a turning of the $+1$ line. Raising e to a real variable power gives a variable number represented geometrically by a variable

length in the real number line. Representing a real number by N , then $\log N$ is the index of e required to make the power equal to N , or to cause a line of length e to stretch to a length N ; before this stretch is made the line e can be raised to the imaginary power $2j\pi$ any number of times, that is, it can be rotated through any number of complete terms without the final result being affected. This is the geometrical meaning of the fact that a real number has an infinite number of logarithms obtained by successive additions of the imaginary logarithm $2j\pi$ to the real logarithm.

e and π .—The reader will see by now how largely the two numbers e and π enter into mathematics. Of these π is generally encountered early in mathematical study as the perimeter of a circle of unit diameter or the area of a circle of unit radius, while, in the conventional textbooks, e appears as a rather mysterious number which is the base of a particular system of logarithms. We have seen that the geometrical significance of e is as clear and simple as that of π ; e is the x ordinate of the rectangular hyperbola $xy=1$ that gives unit area function. We have seen, too, that the curve of $xy=1$, twisted through 45 degrees, becomes $y=\sqrt{(x^2-2)}$, that when x is less numerically than 2 the curve is equivalent to a circle in an imaginary plane, and that π is the measure of the imaginary area function that gives an x ordinate of -1 . Thus the relation of e and π is mutually inverse, e is defined by unit area function or logarithm, π is an imaginary area function or logarithm that gives minus one horizontal ordinate. Symbolically this mutually inverse relation is shown by the equations $\log e=1$, $\log -1=j\pi$.

We have shown how the value of e can be rapidly calculated to a high degree of approximation. The calculation π is not so easy, but we may as well indicate how this calculation can be made independently of geometrical measurement. We saw on p. 206 that the integral of $\frac{1}{1+x^2}$ is $\arctan x$. Further, by reasoning similar to that used to obtain a power series for $\frac{1}{1+x}$ on p. 220, we find that $\frac{1}{1+x^2}=1-x^2+x^4-x^6\ldots$ provided that x is less than 1. Assuming that the integral of

$\frac{1}{1+x^2}$ is the same as the result of integrating the equivalent series term by term, we have

$$\text{arc tan } x = x - \frac{x^3}{3} + \frac{x^5}{5} - \frac{x^7}{7} \dots$$

When $x=0$, the corresponding angle, $\text{arc tan } 0$, is 0 also, so that no constant of integration comes into the equation. This expansion of $\text{arc tan } x$ is by the properties of the geometric series true provided x is less than 1, and it can, as for the $\log(1+x)$ series, be shown to be true when $x=1$. But when $x=1$ or the tangent of an angle is 1, a value of the angle is $\frac{\pi}{4}$ (45 degrees). We thus have

$$\frac{\pi}{4} = 1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7}$$

from which it would be possible to obtain a value of π , although the series converges so slowly that a very large number of terms would have to be taken for even a fair approximation. The derivation of other and more rapidly converging series, which depend upon the arc tan expansion, is explained in text-books on the calculus. Much misdirected and useless effort has been spent in the evaluation of π , and it has been calculated to over 300 places of decimals. The value to 5 decimal places, viz. 3.14159, is sufficiently accurate for really all practical purposes.

It can be proved that both e and π are irrational and transcendental. The discovery of a simple geometrical construction to obtain a straight line equal in length to π is the old classical problem of squaring the circle, which has only quite recently proved to be impossible, but which, for years after this proof, was a favourite occupation for amateur mathematicians. Much interesting information about one of these amateurs who thought he had proved that $\pi=3\frac{1}{8}$ will be found in De Morgan's "Budget of Paradoxes."

Damped Vibrations. Let us return to the exponential function with a real negative index, and consider $x=Ae^{-bt}$. Here the exponent $-bt$ is proportional to time, and the function shows that x has an initial value A when $t=0$ and

$e^{bt} = 1$, and that this value gradually diminishes and approaches zero as t becomes very large. The rate of change of x , that is, the rate at which it decays or dies away to zero, is $\frac{dx}{dt}$ or $-bx$, and this rate depends upon the constant b in the exponent. The differential equation for a dying-away effect of this kind is evidently $\frac{dy}{dx} + bx = 0$.

Now, consider the function $Z = ae^{j\omega t}$, in which the imaginary exponent is j times constant ω , multiplied by time. We know that Z is a when $t=0$ and that for any other time, t , $Z = a (\cos \omega t + j \sin \omega t)$, a complex number of amplitude a , inclined to the real number line at an angle ωt or rather, turned from this line by an amount ωt . $Z = ae^{j\omega t}$ therefore represents geometrically a line of length a rotating about one end with an angular velocity of ω radians a second. It corresponds to the rotating line of Fig. 52, the projections of which on the vertical and horizontal vary periodically according to simple harmonic functions, and as a matter of fact $a \cos \omega t$ is geometrically the horizontal projection, and $a \sin \omega t$ the vertical projection.

Now let us combine these two exponential functions and write $Z = a \times e^{-bt} \times e^{j\omega t}$; ae^{-bt} may be considered as before to be the amplitude of the complex number whose argument or angle is given by $e^{j\omega t}$, but this amplitude is now variable in time; it is gradually diminishing to zero. We have also $Z = ae^{-bt} \cos \omega t + jae^{-bt} \sin \omega t$, which shows that this uniformly rotating but continually shrinking line gives rise to two new kind of functions represented geometrically by its horizontal and vertical projections $ae^{-bt} \cos \omega t$, and $ae^{-bt} \sin \omega t$. This is shown in Fig. 77, the line OP rotating with constant angular velocity ω is continually shrinking because its length is ae^{-bt} . The end of the line moves, not in a circle, as it does in Fig. 52, but in a spiral curve. The angle that the spiral curve makes with the rotating line OP at any instant evidently depends upon the ratio of the speeds with which P is moving tangentially to the spiral and towards the centre of rotation O. The tangential speed is ω multiplied by the length OP, the radial speed is the differential coefficient of ae^{-bt} the length

of OP or $-b$ multiplied by this length. Both speeds are therefore proportional to the length OP , and hence the angle the spiral curve makes with OP is constant. The curve is therefore called an equiangular spiral.

When the line OP is vertical the value of the vertical projection will be a maximum, it will equal e^{-bt_1} where t is the time corresponding to this position. When OP advances

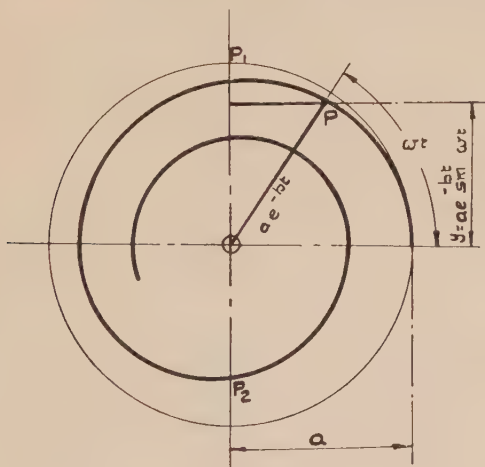


FIG. 77.

to the lower vertical position OP_2 the time will have increased by half the periodic time of the rotation $\frac{T}{2}$ or $\frac{\pi}{\omega}$ and the length

OP_2 will become $e^{-b(t_1 + \frac{\pi}{\omega})}$ or $e^{-bt_1} \times e^{-\frac{b\pi}{\omega}}$. The length of the projection from the position OP_1 to OP_2 will therefore be

decreased by the factor $e^{-\frac{b\pi}{\omega}}$. The next maximum projection on the vertical, upward will be given by a second multiplication by $e^{-\frac{b\pi}{\omega}}$. Hence the extreme swings of the end of its

vertical projection diminish by a constant ratio $e^{-\frac{b\pi}{\omega}}$. This is known as the decrement, and the logarithm of this quantity, $\frac{b\pi}{\omega}$, is called the logarithmic decrement. It should be noted

that it increases as b increases and diminishes as the angular velocity ω increases. The graph of the function $y = ae^{-bt} \sin \omega t$ is shown in Fig. 79. It is a shrinking sinusoidal curve that is bounded by the two lines $y = e^{-bt}$ and $y = -e^{-bt}$. It should be noted that although the amplitude or maximum and minimum values of the graph are continually shrinking, the distances between zero values are constant, as they are in a sinusoidal curve representing simple harmonic motion.

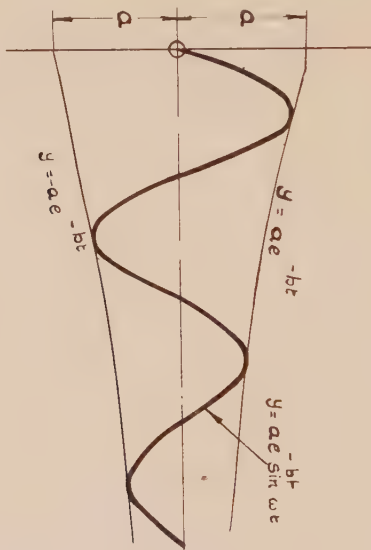


FIG. 78.

Fig. 78 is the graph of what objectively is a damped vibration in which, in addition to a controlling force on the vibrating body proportional to its displacement, there is an additional force proportional to its velocity. This condition obtains very approximately when the bob of a pendulum moves in a viscous liquid. Let us put this dynamical statement into symbolic form: if the mass of the vibrating body is m , then its accelerational force will be $m \times \frac{d^2y}{dt^2}$, this is balanced by the

central force proportional to y , or Fy , plus what is called a damping force proportional to the velocity or $b \times \frac{dy}{dt}$. The differential equation corresponding to the dynamical condition is therefore

$$m \frac{d^2y}{dt^2} + b \frac{dy}{dt} + Fy = 0.$$

Let us assume that this equation is satisfied by $y = Ae^{pt}$, or rather let us try whether this is so. To do this we work out the terms of the differential equation by the rules for differentiation. $\frac{dy}{dt} = py$, so that the $b \frac{dy}{dt}$ term is bpy ; $\frac{d^2y}{dt^2}$ is p^2y , so the $m \frac{d^2y}{dt^2}$ term is mp^2y . Thus if $y = Ae^{pt}$ fits the equation we must have

$$mp^2y + bpy + Fy = 0$$

and as y is not always equal to 0, $mp^2 + bp + F = 0$. Now m , b , and F are constants, so that, as p is an unknown quantity whose value we wish to determine, this last equation is an algebraic quadratic, and the values of p which make it true, will also make $y = Ae^{pt}$ a solution of the differential equation. We can solve the quadratic by means of the general formula given on p. 151. The solution is

$$p = \frac{-b \pm \sqrt{b^2 - 4Fm}}{2m}.$$

Now if $4Fm$ is greater than b^2 , the quantity $\sqrt{b^2 - 4Fm}$ will be negative and we shall have two roots of the quadratic which are conjugate complex numbers like $r + js$ and $r - js$, so that the solution Ae^{pt} will be true for either or both of these values of p , and we can say that our differential equation is satisfied by

$$y = A\{e^{(r+js)t} + e^{(r-js)t}\}$$

or

$$y = Ae^{rt}\{e^{+jst} + e^{-jst}\}$$

But we know that the two terms in the bracket in which e has an imaginary exponent are equivalent to a circular function $2 \cos st$. Hence the condition b^2 is less than Fm , that is, the damping is small, means that the body executes

what is called damped vibrations, of continually diminishing amplitude. The frequency of these damped vibrations depends upon s , that is, upon the quantity $(b^2 - Fm)$ in the solution to the quadratic equation, and this frequency is less than that which would correspond to $b=0$, in which there is no damping and the motion is simple harmonic.

If the quantity $(b^2 - Fm)$ in the solution to the quadratic is positive, that is, if the damping is large, both solutions will be real and the solution to the differential equation will come out like $y = A_1 e^{-rt} + A_2 e^{-st}$, and this will mean that the body will not vibrate but will gradually approach its neutral position, when $y=0$, and the control force disappears. The value of b that just prevents vibration is said to give rise to critical damping.

The foregoing discussion, although only superficial, has taken us into the field of the technique of the solution of differential equations. The important point to grasp is that the character of the motion of the body which we have represented by the differential equation depends upon the values of the mass, control force, and damping constants, as these determine whether the solutions of the quadratic or auxiliary equation are to be real or complex. Only if these solutions are complex do we get vibration of the body about a central position. This differential equation is discussed fully in works on the calculus, and if the reader pursues his studies beyond the scope of this book he will be well advised to spend a good deal of time thereon.

CHAPTER XI

THE APPLICATION OF MATHEMATICS TO STATISTICS

Errors. In the last few chapters we have been considering functional relationships, which, in applied mathematics, are the symbolic statements of those objective uniformities that are called physical laws. Thus the dynamical condition of the simple harmonic motion of a body can be stated as a differential equation, which contains implicitly the law to which this motion is subject, and, by mathematical processes, this implicit law can be discovered and stated explicitly. We have already stated that the final end of the discovery or enunciation of a mathematical function is, in applied science, numerical evaluation. Thus, the arc tan function, expressed as a series, is the basis of the calculation of π whereby its value can be obtained to a standard of accuracy which can transcend that of any actual measurement.

Attention has already been drawn to the approximate character of physical measurements. Whereas we can calculate the value of π to any desired standard of approximation, the accuracy obtainable in its evaluation by actual measurement of the perimeter of a circle would be limited. If a large number of measurements of the same physical quantity are made by a number of observers and in different ways, then although there is one definite true value of the quantity, none of the measured values will be exactly equal to its true value because no measurement can be absolutely accurate. The difference between an actual measure and the true value may be called the error of the measurement. Now errors of measurement may be of two kinds, due to causes which are called assignable and unassignable. Suppose that a measure of length were made by a large number of persons all using the same inaccurate scale. The error of the scale would affect all the measurements equally, because those using the scale would consider it to be accurate. If, however, each person used his own scale, then, although it is possible that each scale would depart from exact accuracy, it is highly improbable that

the scale errors would all be equal in magnitude and sign. When the accuracy of a number of measurements is affected by a single cause operating similarly in respect to each of them, this is known as an assignable cause. When the accuracy is affected by causes so diverse in number and unknown in character that no quantitative significance can be given to them, their resultant is said to be unassignable. A typical example of the operation of unassignable cause arises in the tossing of a coin. Although the result of the toss is believed to be due to the resultant of a number of causes all acting together, these are so unknown that the result is generally said to be a chance one. It is a fundamental postulate of a branch of mathematics called probability that the chances of head and of tail as the result of a toss of a coin are equal, that is to say, that if a number of tosses were made, then, as this number increased without limit, the ratio of the number of heads to that of tails would approach the value of unity, and, as we shall shortly see, when the errors of a number of measurements are subject to what may be called chance rather than to an assignable cause operating alike on all the measurements, it can be inferred that the probability of the occurrence of an error of stipulated magnitude will be functionally related to this magnitude. A very close analogue of errors of measurement is found in the diversity of the physical characteristics of articles that are manufactured all to conform to some standard of uniformity. This standard is not actually attained, and deviations therefrom are governed by conditions and circumstances very like those applicable to errors of measurement.

Averages. Let us think of a set of measurements of the same quantity, all different, because of errors in the measuring process. It is easy to conceive the idea of a single number representative of the whole set. The most obvious is the arithmetical average, commonly called the mean. If the various measurements of the set are N in number, and are indicated by x_1, x_2, \dots, x_n , then the mean is equal to the sum of all the x 's divided by N , and it is denoted by \bar{x} . Thus $\bar{x} = \frac{\Sigma x}{N}$; where Σx means the sum of all the x 's.

An alternative to this is the geometric mean, it is the N th root of the product of all the x 's. It is easy to see that the logarithm of the geometric mean is the arithmetic mean of the logarithms of the x 's. The arithmetic is greater than the geometric mean. It is easy to show that this is so for two quantities, x , and $x+h$, the mean of which is $x + \frac{h}{2}$. For the square of the mean is $x^2 + xh + \frac{h^2}{4}$ and this is greater than the square of the geometric mean, $x(x+h) = x^2 + xh$.

Another average value is the median which is assessed in quite a different way. Suppose that a set of magnitudes is arranged, or arrayed as it is sometimes termed, in order of magnitude, then the middle magnitude of the array is the median average.

Deviation. The quantity $x - \bar{x}$, or the difference between any magnitude of a set, and the mean of the set is called the deviation of x , and may be denoted by d . d is positive if x is greater than \bar{x} , and negative if it is less. The algebraic sum of all the d 's is evidently zero, for, if there are N of them, this sum is $\Sigma x - N\bar{x}$ and $N\bar{x}$ is by definition equal to Σx .

An important secondary characteristic of a set of quantities, all normally equal, is an average of all the deviations. A mean value in the ordinary sense would be valueless because the deviations add up to zero. We might obtain a mean value by ignoring sign and treating all the d 's as positive. It has been found more useful, however, to average the squares of the deviations; this destroys the differences of signs because all squares are positive. The square root of the mean of the squares of the deviations is called the standard deviation,

denoted by σ . Thus $\sigma^2 = \frac{\Sigma d^2}{N}$. σ evidently is a measure of the extent of the spread of the quantities about the mean value.

An important algebraic property of the standard deviation should be noted. Let X be an arbitrarily fixed value from which a kind of false deviations are measured. If $X - \bar{x} = K$, then the false deviation corresponding to d will be $d + K$, the square of this false deviation will be $d^2 + 2dK + K^2$, and the sum of all these squares will be $\Sigma d^2 + 2(\Sigma d)K + NK^2$. Now, as K

is a constant for all d 's, it is a constant factor in ΣdK , which therefore equals $K\Sigma d$. But as Σd , or the algebraic sum of all the deviations, is zero ΣKd is zero also, so we have $\Sigma d^2 + NK^2$ = sum of squares of false deviations. It therefore follows that if deviations are taken from a false mean, itself having a deviation of K , then the mean square of the false standard deviation is equal to $\sigma^2 + K^2$. This is important in actual calculation. Thus, if a set of measurements have a mean value of 21.07 it is quicker to calculate false deviations from a false mean of 20, and correct the false standard deviation by subtracting $(1.07)^2$, than to reckon true deviations from the actual mean, 21.07.

Frequency Distribution. Although the mean and the standard deviation of a set of observed quantities are important, they do not tell us very much. Many more implications of the set are brought out by means of a graph of what is called cumulative frequency distribution. This graph is plotted in the following way. Let horizontal ordinates represent observed quantity x and the corresponding vertical ordinate the number n of the quantities out of the total, N , that are less than x . Fig. 79 (a) is a graph of this kind. The vertical ordinates extend from 0 to N , and a horizontal line through the point on the vertical axis, $\frac{N}{2}$, defines an x value evidently corresponding to the median. Let us consider the differential properties of the graph. We have n_1 quantities are less than x_1 and n_2 less than x_2 . The difference between n_2 and n_1 or δn is therefore the number of quantities not less than x_1 but less than x_2 , or the number lying in the interval $\delta x = x_2 - x_1$. If the base of the graph is divided into a large number of equal parts δx , then to each part will correspond a number, δn , of the quantities that lie in the interval δx , which is defined by its central x value. This number is called the frequency corresponding to the x value, and graphically it is measured by the slope of the graph, the greater the slope the greater will be the number of observed quantities lying in a constant interval δx . The x value that defines the point of greatest slope, M_o in the figure, is called the modal average or the mode. It is the x

value in the neighbourhood of which the frequency is greatest.

Let us suppose that the mean, \bar{x} , of all the observed quantities is defined by the point Me on the graph. Then δn or

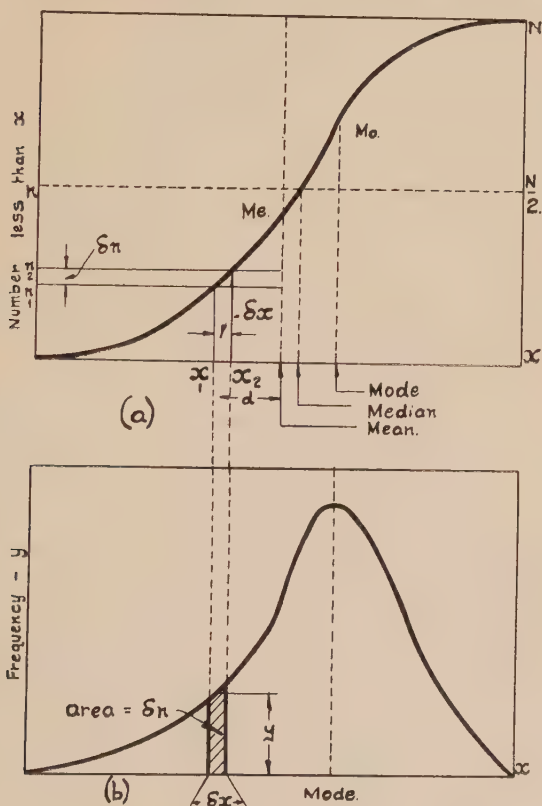


FIG. 79.

$n_2 - n_1$ will be the number of quantities that have an average deviation d , which is defined by the x of the δx corresponding. d is here shown negative, and it is clear that the sum of all the negative d 's is the quantity $\sum d \times \delta n$, is the area contained by the graph to the left of Me , the vertical through Me and the horizontal through $n=0$. Similarly, the sum of all the

positive deviations is the area contained by this vertical, the graph to the right of it, and the horizontal through $n = N$. But by the property of the mean, the sums of positive and negative deviations are equal, hence the vertical through Me makes these two areas equal. It is the average width of the graph from the vertical axis.

If we draw a graph of the slope function of Fig. 79 (a) we get a derived graph like Fig. 79 (b), which is called a frequency curve. The point of maximum frequency defines the mode. Any ordinate y of the frequency curve represents to scale the slope of the cumulative frequency curve for the x value corresponding. But this slope $\frac{dn}{dx}$ is practically the same as $\frac{\delta n}{\delta x}$ where δn and δx refer to small but measurable increments of the ordinates of the cumulative curve. As δx of the one curve may be considered to be the same as that of the other, and $y = \frac{\delta n}{\delta x}$, for the frequency curve it follows that δn , the number of quantities corresponding to the interval is given by $y \times \delta x$ which is geometrically represented by a vertical strip of the frequency curve of width δx . Hence the sum of all such strips, or the total area under the frequency curve represents, to scale, the total number of quantities, N , to which the curve refers.

Normal Frequency Curve. Let us now consider the matter of errors of measurement in a general sort of way. If a very large number of measurements of the same kind of quantity are made under similar conditions, both experience and intuition seem to show that the errors will be subject to a kind of law ; in the long run, a positive error of given amount will be as likely as a negative error of the same amount, and large errors will be less likely than small ones. If these assumptions or postulates are granted it follows that the frequency curve will be symmetrical about its maximum ordinate, which will define the three averages, mean, median, and mode, all equal. Further, if the horizontal ordinates are measured as deviations, in other words, if the vertical axis of co-ordinates is the ordinate through the mean, then the curve will be represented symbolically by a function of the deviation, which has a maximum value when this is zero, and which dies away

towards zero as this deviation increases in either a positive or negative direction.

We can see, without much difficulty, that a function satisfying this condition is $y = Ae^{-\frac{x^2}{h^2}}$ where x stands for the deviation or error. Whether x is positive or negative x^2 is always positive, the exponent of e is always negative, so that as x increases from 0 in either direction, the negative power of e increases. This, of course, is not the only function that will satisfy the fundamental postulates regarding frequency of or what is practically the same thing, probabilities of errors. This particular function, however, does rest upon proofs too complicated to be given here, which themselves are based upon even simpler, but still unprovable, postulates. Further, the function has for a long time been used as the basis of practical calculations and has been found to lead to useful results. Hence, although we cannot say with certainty that $y = Ae^{-\frac{x^2}{h^2}}$ represents the normal frequency of errors, the mathematical statement can be considered to be supported by some measure of theoretical proof and also by the pragmatic sanction of successful practical use.

The graph of the function $y = Ae^{-\frac{x^2}{h^2}}$ is shown in Fig. 80. It is known as the normal curve of error, and as the Gaussian frequency curve, because it was first used by the mathematician Gauss. The curve is symmetrical with respect to the vertical axis of zero error or deviation, and it extends indefinitely in both directions, being asymptotic to the horizontal axis. The curve itself depends upon the value of the constant h in the exponent, for fixing a value of y , that is, of the exponent, the value of the x required to produce this y must increase proportionally to h . Thus, as h is increased, the width of the curve at any height increases correspondingly as shown by the dotted graph. All the curves comprised by the functional formula with the parameter h , taking any value, are therefore similar in shape. h therefore defines the width of the spread of errors about the zero datum; it is a kind of indication of the accuracy of the method of measurement. We shall see its precise significance very soon.

The area comprised between the curve and the horizontal axis, extending indefinitely in each direction, tends to a definite limit, and it can be shown by a most elegant application of the integral calculus, too difficult to be given here, that when $A=1$, this area is $\frac{1}{2}h \times \sqrt{\pi}$.

Let us investigate the differential properties of the normal frequency curve. We can calculate $\frac{dy}{dx}$ by the rule for differentiation of a function of a function, and we find that

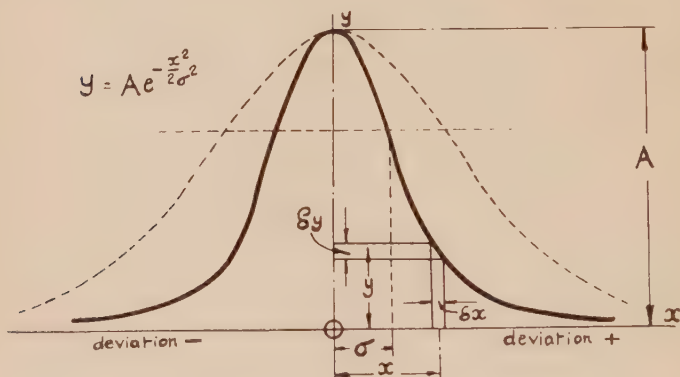


FIG. 80.

$\frac{dy}{dx} = Ae^{-\frac{x^2}{h^2}} \times -\frac{2x}{h^2} = -\frac{2xy}{h^2}$. Since $\frac{dy}{dx}$ is very approximately $\frac{\delta y}{\delta x}$ where δy and δx refer to measurable increments, it follows that $\frac{\delta y}{\delta x} = -\frac{2xy}{h^2}$.

We have seen that $y\delta x$, representing a vertical strip of a frequency curve, is a measure of the number of quantities corresponding to the interval δx , and having an average deviation or error x , corresponding to the y . Thus, assigning an average deviation x , the sum of the squares of the deviations having this average value will be $x^2 \times$ the number $y\delta x$ or $x^2 y\delta x$. But, as $\frac{\delta y}{\delta x} = -\frac{2xy}{h^2}$, $xy = -\frac{\delta y}{\delta x} \times \frac{h^2}{2}$, so that $x^2 y\delta x = -\frac{h^2}{2} \times x\delta y$. But the sum of all such quantities as $x^2 y\delta x$ is the sum of the

squares of all the deviations and is equal to $\sigma^2 N$ where σ is the standard deviation, and as positive and negative deviations are distributed in exactly the same way we can say that $\sigma^2 \times (\text{area } \frac{1}{2} \text{ curve})$ is equal to the sum of all such quantities as $x^2 y \delta x$ for positive deviations. But this sum is also that of such terms as $-\frac{h^2}{2} \times x \delta y$. Now $x \delta y$ for positive values of x is a horizontal strip of the half curve shown in the figure, with a negative sign because δy is negative or y is decreasing. Hence the sum of all such quantities as $-\frac{h^2}{2} \times x \delta y$ is equal to $\frac{h^2}{2} \times (\text{area } \frac{1}{2} \text{ curve})$. It follows therefore that $\sigma^2 = \frac{h^2}{2}$ and $h^2 = 2\sigma^2$.

The normal frequency curve can therefore be written as $y = Ae^{-\frac{x^2}{2\sigma^2}}$ where σ is the standard deviation.

The $\frac{dy}{dx}$ of the normal frequency curve now becomes $-\frac{xy}{\sigma^2}$

The second differential coefficient $\frac{d^2y}{dx^2}$ can be found by dif-

ferentiating this product; it is $-\frac{1}{\sigma^2} \left(-x \times \frac{xy}{\sigma^2} + y \right)$. Now the slope of the curve will be a maximum when this second differential coefficient is zero, or when $-\frac{y}{\sigma^2} \left(-\frac{x^2}{\sigma^2} + 1 \right) = 0$.

As y will evidently not be zero for maximum slope we have $-\frac{x^2}{\sigma^2} + 1 = 0$ or $x = \sigma$, so that σ the standard deviation defines the points of maximum slope on the curve. At this point

$y = Ae^{-\frac{\sigma^2}{2\sigma^2}} = Ae^{-\frac{1}{2}}$, this y is the same for all curves that have the same maximum value A , and, as these curves are all similar in shape it follows that vertical ordinates through the two points of maximum slope, the one positive and the other negative, corresponding to $x = +\sigma$ and $x = -\sigma$, contain between them a definite fraction of the total area of the curve which can be shown to have the value 0.6745. This means that if a large number of deviations correspond in their frequency to

the normal curve, then 0.6745 of them do not exceed σ , the standard deviation.

The practical consequence of this last statement is very important. Suppose that a set of observed quantities are plotted to give a cumulative frequency curve, like Fig. 81. If the actual frequencies conform to the normal distribution then the horizontal line through $n = 0.5N$ defines the mean, and two horizontal lines through $n = 0.163N$ and $n = 0.837N$ comprise between them a fraction equal to 0.674N.

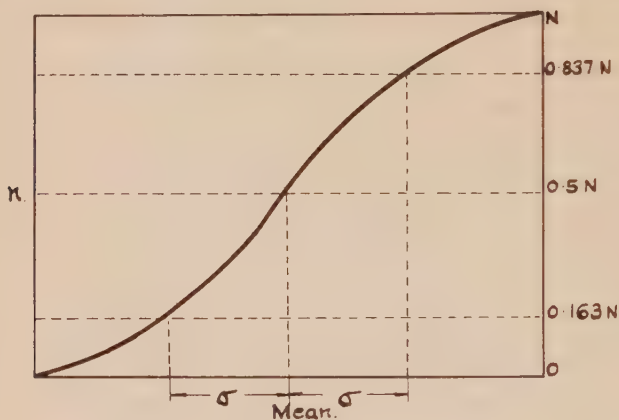


FIG. 81.

These two lines, therefore, define points on the curve for which the x values have a deviation from the mean of σ and $+\sigma$. Thus the standard deviation is indicated directly provided the frequency distribution follows the normal law. In practice this can be tested very simply by the use of a special kind of graph paper with the vertical divisions so marked that any normal frequency distribution gives a straight line cumulative frequency graph. Provided the plotted graph is straight, the standard deviation can be directly read off, and evidently it is indicated graphically by the inclination of the graph relative to the vertical.

Principle of Least Squares. We may conclude this chapter by referring briefly to an important corollary of the normal

law of error which is known as the Principle of Least Squares. Suppose that a number of observations, say 3, are taken of the same quantity under the same conditions, so that the constant σ in the error function is the same for them all. Let the unknown errors of these measurements be called a_1 , a_2 , and a_3 . The likelihood or probability of the occurrence of the

error a_1 is proportional to $e^{-\frac{a_1^2}{2\sigma^2}}$, and so with the probabilities of the errors a_2 and a_3 . The probability of their all occurring in a set of three measurements is less than either component probability; the combined probability is equal to the product of the components, that is, to e raised to a power $-\frac{1}{2\sigma^2}(a_1^2 + a_2^2 + a_3^2)$ because in multiplying powers of e we add

the indices or exponents. Now the smaller a negative exponent the greater the power will be, so the greatest probability of the occurrence of the errors is given by the condition that $a_1^2 + a_2^2 + a_3^2$ is a minimum. It follows that the most probable value of the quantity measured is that which gives the squares of the assumed errors referred to it or the deviations from it a minimum sum. But we have seen that the sum of the squares of the deviations from any assigned false average is equal to $\Sigma d^2 + NK^2$ where d stands for a deviation from the mean value. N the number of observations and K the difference between the false average and the true one. As K^2 is always positive, it follows also that the sum of the squares of the deviations from any datum or false average is the least when $K=0$ and when the datum is not false but identical with the mean value. Hence the mean of a set of observations is the most probable value of the quantity measured.

Again, suppose that two measurable variables are known to be connected by a functional relationship $y=a+bx$, and suppose that three sets of values (x_1, y_1) (x_2, y_2) (x_3, y_3) have been actually measured. Now if these sets of values are plotted into points on graph paper, then, generally, the three points will not all lie in a single line owing to errors of measurement. It will be possible to draw three straight lines, each through a pair of points, and neither of these lines will represent the actual function. We can use the principle of Least Squares

to discover, from the measured values, the best values of a and b in the equation $y = a + bx$. For, taking a measured value x_1 , the true y corresponding to it is $a + bx_1$, and the error of the measured y_1 is $(y_1 - a - bx_1)$. Similarly for the other two observed y values.

The most probable values of the unknown constants a and b are those which make the sum of the squares of the y errors a minimum. The y_1 error is $(y_1 - a - bx_1)$. We have not yet worked out the square of an algebraic expression with three terms, or a trinomial, as it is called, but it is not difficult to see that this square is $y_1^2 + a^2 + b^2x_1^2 - 2ay_1 + 2abx_1 - 2bx_1y_1$. We note that in this expression that x_1 and y_1 are numbers. The squares of the other errors will be similar expressions containing the same a and b but different y 's and x 's. The sum of the squares will therefore be

$$\Sigma y_1 + Na^2 + b^2 \Sigma x^2 - 2a \Sigma y + 2ab \Sigma x - 2b \Sigma xy$$

where N is the number of sets of values, three in the case we supposed.

How are we to find the values of a and b that will make this quantity a minimum? Its value will be affected both by changes of a and b . Suppose b remains constant. One condition for a minimum will be that the differential coefficient with respect to a is equal to zero. We carry out this differentiation by the usual rules, but we treat b , as well as the x 's and y 's which are known numbers, as constants. The result of this is $2Na - 2\Sigma y + 2b\Sigma x$, and equating this to zero we have for the first condition $Na - \Sigma y + b\Sigma x = 0$. Differentiating with respect to b and treating a as constant we get $2b\Sigma x^2 + 2a\Sigma x - 2\Sigma xy$, and equating this to zero, $b\Sigma x^2 + a\Sigma x - \Sigma xy = 0$. These two equations to zero are a pair of simultaneous equations with a and b as the unknowns, and all the other terms actual observed numbers. Hence the equations can be solved as explained in Chapter VI, and the most probable values of a and b calculated. This is much more accurate than the corresponding geometrical process of drawing a line through a number of plotted points by judgment of the eye.

INDEX

- Abscissa, 107
- Absolute term, 141
- Acceleration of S.H.M., 179
- Accuracy of measurement, 60
- Acute angle, 32
- Addition, 10
 - of fractions, 55
 - of vectors, 51
- Algebraic functions, 91
 - symbolism, 26
- Amplitude, 160
- Angle, 32, 78
 - , measure of, 46, 159
- Angular velocity, 159
- Approximate solution of equation, 148, 151
 - square roots, 67
- Approximation, 59
- Arbitrary constant, 203, 209
 - function, 210
- arc sin, 50
- Area, 36
 - by integration, 207
 - function, 122
 - — of hyperbola, 213
- Argument of complex number, 81
- Arithmetic, 8
 - mean, 238
 - series, 99
- Array, 239
- Assignable causes of error, 237
- Associative law, 11
- Asymptote, 116
- Auxiliary equation, 236
- Average, 238
 - height of graph, 208
 - value of periodic function, 172, 182
- Axiom, 30
- Beats, 171
- Binary system, 10
- Binomial, 22
 - theorem, 23
- Bolyai, 34
- Brackets, 12
- Calculus, differential, 185
 - , integral, 202
 - operational, 194
- Cardan's method, 151
- Cardinal numbers, 8
- Cartesian notation for vectors, 53
 - co-ordinates, 107
- Casting out nines, 19
- Centre of area, 208
- Centroid, 208
- Chance, 238
- Characteristic of logarithm, 73
- Cis θ , 81
- Circle, 37, 107
 - , division of, 47
- Circular functions, 180
 - measure, 48
- Circulating decimals, 59, 102
- Coefficient, 22
- Commensurable ratio, 40
- Commutative law, 11
- Common logarithm, 68, 220
 - parabola, 114
- Completing the square, 150
- Complex numbers, 80
- Composition of S.H.M.'s, 107
- Complementary function, 211
- Congruent triangles, 35
- Conjugate complex numbers, 81
- Consistent equations, 139
- Constant, 86
 - of integration, 203
- Continued fractions, 40, 62, 66
- Continuity, 126
- Continuous function, 93
- Convergents, 62
- Convergence of series, 101

- Cosine, 44
 - of multiple angles, 82
- Cosecant, 45
- Cosh u , 223
- Cotangent, 145
- Counting, 8
- Critical damping, 236
- Cube, 20
- Cubic equation, 151
- Cumulative frequency curve, 240
- Curve tracing, 127

- δ , 60
- Damped vibrations, 231
- Decimal fractions, 57
 - system, 9
- Definite integral, 204
- Definitions, 30
- De Morgan, 231
- Dependent variable, 85
- Derived function, 121, 184
 - unit, 89
- Determinant, 138
- Deviation, 239
- Differential coefficient, 185
 - equation, 191, 209
- Differentiation, 186
- Digit, 8
- Dimensions, 89
- Directed numbers, 69
- Discontinuity, 93
- Divergence of series, 101
- Divisibility of numbers, 18
- Division, 15
 - of approximate numbers, 62
 - of fractions, 156
- Duodecimal system, 9
- Dynamics of S.H.M., 180

- e , 97, 217, 230
- Elimination, 137, 210
- Ellipse, 128, 167
- Empirical graphs, 129
- Equal roots, 145
- Equations, 87, 113
 - , consistent, 139
 - , cubic, 151
 - , integral, 140
 - , quadratic, 152
 - , simple, 135
 - , simultaneous, 136
- Equations, transcendental, 134
- Equiangular spiral, 233
- Equivalent triangles, 35
- Error curve, 242
- Errors of observation, 130, 237
- Euclid's Elements, 18, 30
 - parallel postulate, 33
- Evaluation of series, 104
- Even function, 113
- Evolution, 25, 74
 - of complex numbers, 82
- $\exp x$, 91
- Explicit function, 86
- Exponential function, 97, 190, 218
 - theorem, 219
 - value of $\sin \theta$, 227
- Extrapolation, 130

- Factorial, 17
- Factors, 17, 141
- First moment, 208
- Fitting graph, 130, 248
- Forced vibrations, 212
- Formal solution of equation, 133
- Formula, 87
 - for solving quadratic, 151
- Fourier's theorem, 180
- Fractional errors, 60
 - index, 67
- Fractions, 55
 - , continued, 40
 - , decimal, 157
 - , vulgar, 57
- Frequency of S.H.M., 157
 - errors, 240
 - curve, 241
 - distribution, 240
- Function, 85
 - , algebraic, 91
 - , arbitrary, 211
 - , area, 122
 - , circular, 158, 180
 - , complementary, 211
 - , continuous, 92
 - , derived, 121
 - , discontinuous, 93
 - , even, 113
 - , exponential, 97, 218
 - , explicit, 86
 - , hyperbolic, 221
 - , implicit, 86
 - , integral, 92

Function, linear, 109
 —, logarithmic, 220
 —, many-valued, 183
 —, odd, 113
 —, periodic, 156
 —, rational, 92
 —, single-valued, 92
 —, slope, 121
 —, transcendental, 96
 —, trigonometrical, 158
 Fundamental, 180

Gauss, 47
 Gaussian frequency curve, 243
 Geometry, 30
 Geometric mean, 239
 — series, 99
 Graph, 106
 — fitting, 130, 248
 Graphical solution of equations, 146

Harmonic analysis, 181
 — function, 158
 — series, 101
 Harmonics, 180
 Homogeneous function, 88
 Horner's method, 154
 Hyperbola, 116, 213
 Hyperbolic functions, 221
 — logarithm, 220
 Hypotenuse, 41

i , 76
 Identity, 134, 141
 Imaginary number, 80
 — index, 166, 227
 — logarithm, 229
 Implicit function, 86
 Incommensurable numbers, 57
 — ratios, 40
 Independent variable, 85
 Index, 20
 —, fractional, 67
 —, irrational, 96
 —, negative, 72
 —, notation, 24
 Infinite solutions, 137
 — series, 100, 193
 Infinity, 17

Inflecting point, 125
 Integral calculus, 202
 — function, 92
 Integration, 204
 — of differential equation, 210
 Intercept form, 111
 Interpolation, 130
 Inverse proportion, 93, 116
 — circular functions, 50
 — hyperbolic functions, 226
 Irrational numbers, 63

j , 76, 79

Lag, 170
 Lead, 170
 Least squares, 247
 Limiting value, 92
 Linear function, 109
 Lobatschewsky, 34
 Locus, 107
 Logarithm, 24, 68, 73, 220
 Logarithmic decrement, 283
 — function, 97, 190, 213
 — series, 221

Mantissa, 73
 Many-valued function, 183
 Maximum values, 125, 195
 Mean, 238
 Median, 237
 Minimum values, 125, 195, 248
 Mode, 240
 Modulus of complex number, 81
 Moment of inertia, 209
 Multiplication, 13
 — of approximate numbers, 61
 — of complex numbers, 78
 — of fractions, 55
 — of irrational numbers, 64

Naperian logarithm, 220
 Natural numbers, 8
 Negative indices, 72
 — numbers, 68
 Newton's method of solving equations, 148
 Non-Euclidean geometry, 33, 36
 Normal frequency curve, 242

Notation of calculus, 186
 Numbers, cardinal, 8
 —, conjugate, 81
 —, complex, 80
 —, imaginary, 80
 —, incommensurable, 48
 —, irrational, 63
 —, natural, 8
 —, negative, 68
 —, positive, 68
 —, puzzles on, 26
 —, rational, 63
 —, real, 68, 72
 Numeration, 8
 Numerical solution of equations,
 133

Obtuse angle, 32
 Odd function, 113
 Order of contact, 147
 — of differential equation, 210
 — of number, 23
 Ordinate, 107

π , 37, 63, 230
 Parabola, 114, 116
 Parallel lines, 26
 — postulate, 37
 Parameter, 86
 Partial fraction, 16
 — product, 16
 Particular integral, 211
 Percentage, 59
 Periodic function, 156
 — time, 157
 Phase, 160, 169
 Plotting graphs, 127
 Point, 31
 — of inflection, 125
 Polar co-ordinates, 107
 — representation of vector, 53
 Positive numbers, 68
 Postulate, 30, 64
 Prime, 17
 Principal value, 183
 Probability, 238, 247
 Progression, 98
 Proportion, 40
 Pythagoras's theorem, 41
 Pythagorean triplets, 42

Quadratic equation, 148
 — irrational, 65
 Quarternious, 84
 Quartic equation, 152

Radian, 48
 Radius of gyration, 209
 — vector, 167
 Rational function, 92
 — number, 63
 Rationalising factor, 65
 Real number, 68, 72
 Reciprocal, 92
 Rectangular hyperbola, 96, 213
 Recurring continued fractions, 66
 — decimals, 59, 102
 Reimann, 34
 Remainder, 16, 141
 — theorem, 141
 Resonance, 212
 Right-angle, 32
 Right-angled triangle, 41
 Root, 25, 66
 — of equation, 83
 — of unity, 133
 Rule of signs, 71

Σ , 202
 Secant, 45
 Second moment, 208
 Secular change, 156
 Sense of vector, 51
 Series, 97
 S.H.M., 158
 Similar triangles, 35
 Simple equations, 135
 — harmonic motion, 158
 Simultaneous equations, 136
 Sine, 44
 Sinh u , 223
 Sinusoidal variation, 160
 Slope, 109, 120
 — function, 121
 Small angles, 161
 Solution of equations, 133, 141
 — of differential equations, 210
 Space, 30
 Spiral, 112, 223
 Square, 14
 — number, 25
 — root, 65

Squaring the circle, 231
Standard deviation, 239, 245
Straight line, 31
— — graph, 109
Subtraction, 11
— of fractions, 55
Successive differentiation, 192
Sum to infinity, 100, 103
Surd, 65
Surface, 31
Symbolic representation of number, 74

Tangent, 121
 $\tanh u$, 223
Test for convergence, 104
Transcendental function, 96
— number, 68
Transversal, 32
Triangle, 35
—, right-angled, 41
Trigonometrical ratios, 42
Turning values, 195

Unassignable causes of error, 238
Undetermined coefficients, 144
Unit, 39, 89
—, derived, 90

Variable, 85
Vector, 50
— addition, 51
— multiplication, 78
— representation of S.H.M., 168
Velocity of wave, 158
— of S.H.M., 176
Vulgar fraction, 57

Wave form, 157
— length, 158
— motion, 158, 210

Zero, 9, 11, 15, 21

W. VA. INSTITUTE OF TECHNOLOGY



* 1800031620 *

QA
37
S88

36890

THE LIBRARY
WEST VIRGINIA
INSTITUTE OF TECHNOLOGY
MONTGOMERY, WEST VIRGINIA 25136

LIBRARY BUREAU CAT. NO. 1169.6

